

VILNIAUS UNIVERSITETAS  
MATEMATIKOS IR INFORMATIKOS FAKULTETAS  
KOMPIUTERIJOS KATEDRA

Baigiamasis magistro darbas  
**Emocijų klasifikavimo metodų tyrimas teksto analizės uždaviniuose**

Atliko:  
Kristina Rasa Tamolė parašas  
Vadovas:  
dr. Linas Bukauskas

Vilnius  
2017

## Turinys

Turinys.....	2
Sutartinis terminų žodynas.....	3
Santrauka.....	4
Summary.....	5
Įvadas.....	6
1. Analitinė dalis.....	9
1.1. Sentimentinės analizės taikymai teksto gavybos uždaviniuose.....	9
1.2. Emocijų analizė – sentimentinės analizės kryptis.....	10
1.2.1. Emocijų analizės raida.....	10
1.2.2. Emocijų aptikimo tekste problematika.....	11
1.2.3. Emocijų vertinimo taksonomija.....	12
1.2.4. Emocijų raiškos internetinėje erdvėje būdai.....	13
1.3. Sentimentų žodynai.....	14
1.3. 1. Sentimentų žodynų apžvalga.....	14
1.3. 2. Sentimentų žodynų generavimo metodai.....	15
1.3. Sentimentinės analizės metodų apžvalga.....	16
1.3.1. Klasifikavimui naudingų požymių atranka.....	17
1.3.2. Mašininio mokymo metodų apžvalga.....	18
1.3.2.1. Prižiūrimasis mokymas.....	18
1.3.2.2. Iš dalies prižiūrimas ir neprižiūrimas mokymas.....	19
1.3.3. Sentimentų žodynu grindžiamų metodų apžvalga.....	19
2. Emocijų analizės eksperimentas.....	21
2.1. Duomenų gavyba.....	21
2.2. Teksto emocijų klasifikavimas.....	25
2.3.1. Semantinio metodo taikymas emocijų nustatymo uždavinyje.....	25
2.3.1.1. Pirminis duomenų paruošimas.....	26
2.3.1.2. Klasifikavimui naudingų požymių vektorius.....	26
2.3.1.3. Lingvistinių taisyklių modeliavimas.....	27
2.3.1.4. Klasifikavimo rezultatai.....	29
2.3.1.5. Ironijos nustatymo algoritmas.....	30
2.3.2. Mašininio-statistinio metodo taikymas emocijų nustatymo uždavinyje.....	31
2.3.2.1. Naudingų klasifikavimui požymių atrinkimas taikant statistinį metodą.....	31
2.3.2.2. Teksto emocijų klasifikavimas taikant regresijos metodą.....	33
2.3.3. Eksperimento rezultatų vertinimo kriterijai.....	35
Apibendrinimas.....	38
Pagrindiniai rezultatai.....	39
Ateities tyrimų gairės.....	41
Literatūros šaltiniai.....	42

## Sutartinis terminų žodynas

1. Sentimentinė analizė (angl. *sentiment analysis*) – jausmų ir nuomonės analizė.
2. BoW (angl. *bag-of-words*) - žodžių krepšelio metodas.
3. TF (angl. *term frequency*) – termino pasikartojimo dažnumas.
4. Duomenų gavyba (angl. *data mining*) – tai žinios, kurias galima gauti iš jau turimų duomenų ir jas atitinkamai apdoroti.
5. Žiniatinklio gavyba (angl. *web mining*) – duomenų gavybos metodų pritaikymas Žiniatinklio sistemoms.
6. Nuomonės gavyba (angl. *opinion mining*) yra kitaip žinoma, kaip sentimentine analizė.
7. Teksto gavyba (angl. *text mining*) - tai duomenų gavybos pritaikymas nestruktūrizuotiems ar pusiau struktūrizuotiems duomenimis.
8. Reguliarios išraiškos (angl. *regular expression arba regex*) – tai taisyklių rinkiniai, aprašantys tekstinį šabloną, pagal kurį randamas reikiamas tekstas arba jo fragmentas.
9. VDU KLC – Vytauto Didžiojo universiteto kompiuterinės lingvistikos centras
10. Taksonomija – matematiškai apibrėžiama kaip klasifikacija, kurioje visi objektai yra išdėstomi į tam tikrą medžio struktūrą. Šio medžio viršūnę paprastai sudaro vienintelis klasifikacijos elementas, kuriam priklauso visi kiti šios klasifikacijos objektai. Taksonas, einantis žemiau už viršūnę, yra labiau specifinis ir išskiria tam tikrą klasifikuojamų objektų aibę pagal tam tikrą požymį.
11. Konotacija (lot. *con-* – „su-“, „kartu“, *notatio* – „pastaba“) – kalbos vieneto (žodžio) šalutinė reikšmė, rodanti kalbos vartotojo santykį su sąvokomis ir jomis išreiškiamais objektais.
12. Žvalgomasis tyrimas – preliminarus tyrimas, kuriame pagrindinis dėmesys sutelkiamas idėjų generavimui ir informacijos, reikalingos problemos/algorithmo formulavimui, paieškai. Atliekamas tuomet, kai tyrimo problema nėra pakankamai aiški.
13. Apriūrinis – išankstinis, nepriklausantis nuo patyrimo.
14. Retorinės figūros – tai kalbos ekspresyvumą didinančios stilistinės priemonės, kurioms būdingas žodžių vartojimas perkeltine prasme, aiškinant vieną reiškinį kitu

## **Santrauka**

Magistro baigiamajame darbe nagrinėjamas lietuviškų žiniatinklio dokumentų emocijų klasifikavimo uždavinys. Šis uždavinys sprendžiamas taikant natūralios kalbos ir mašininio mokymo technologijas. Pagrindinis darbo tikslas – teksto emocijų klasifikavimo metodų tyrimas įgyvendintas remiantis semantinės krypties ir statistiniu metodais, pasitelkiant žodžių bei frazių taksonomiją, lingvistinėmis taisyklėmis grindžiamą modelį, statistiniu metodu nustatytų reikšmingų klasifikavimui požymių sąrašą, kuriems regresijos ir gradientinio nusileidimo algoritmais apskaičiuoti koeficientai. Eksperimento metu nustatyta, kad aukštesnė teksto emocijos klasifikavimo kokybė stebima taikant semantinės krypties metodą, kurio įgyvendinimas grindžiamas emocijų žodyno ir lingvistinių taisyklių modelio taikymu.

## **Summary**

### **Research of Emotion Detection Methods for Text Mining**

This paper investigates the problem of emotion classification in online texts. Lexicon-based model and statistical approach based on machine learning have been developed to classify an online text into one of multiple emotion categories (i.e. “anger”, “disgust”, “fear”, “happiness”, “sadness”, “surprise”).

Leveraging Ekman emotion framework and the RSS feed data of real-world news portals "Delfi" and "Lietuvos rytas", a lexicon-based framework has been developed to associate most common online text terms and affect words with a distribution on a series of emotions. The framework is based on manually designed emotion lexicon and sequential linguistic rules.

Using an annotated set of debate forum posts and news articles, statistical emotion detection model has been designed by extracting patterns that are highly correlated with emotion expression. The statistical approach is based on chi-square and regression algorithms.

The performance of machine-statistical and lexicon-based classifiers has been evaluated and compared. The experimental analysis on the task of emotion classification validates the effectiveness of the proposed lexicon-based model for documents containing traditional affective vocabulary. The proposed affective model can be applied to both – the tasks of classifying emotions and tasks of generating social emotion lexicons.

## **Įvadas**

Didėjant žiniatinklio prieinamumui internetinėje erdvėje pateikiama informacija tampa svarbiu veiksniumi ne tik informuojant apie įvykius, bet ir perteikiant nuomonę, požiūrį, emocijas. Naujienų portaluose išsakoma autorių ir skaitytojų nuomonė, pasižyminti subjektyviu požiūriu, formuoja emocinę internetinės erdvės atmosferą.

Kadangi internetinėje erdvėje perteikiamos emocijos įtakoja žmogaus nuotaiką, psichinę sveikatą ir asmenybės vystymąsi[9], automatizuotas emocijų nustatymas internetiniame tekste tampa aktualiu teksto analizės uždaviniu.

**Tyrimo objektas** – emocijų nustatymo internetiniame tekste metodai ir jų pritaikymas teksto emocijų analizei.

**Mokslinis tyrimo naujumas ir vertė.** Emocijų analizė – vienas sentimentinės analizės aspektų, kuriuo siekiama nustatyti emocijas klase tekste. Emocijų nustatymas yra aktualus marketingo, medicinos, politikos srityse, tačiau rankiniu būdu nustatyti internetinio teksto emociją užima daug laiko ir dėl informacijos gausos yra praktiškai neįgyvendinamas uždavinys. Siekiant nustatyti emociją elektroniniame tekste kuriami automatizuoti teksto emocijų klasifikatoriai. Tam būtinas kokybiškas emocijas indikuojančių žodžių, kitaip vadinamų emociniais terminais, žodynas. Anglų kalba yra sukurta keletas išsamių žodynų, skirtų teksto emocijų klasifikavimo uždaviniui spręsti, tačiau kitomis kalbomis, tame tarpe ir lietuvių, tai vis dar neišspręsta problema, todėl siektinas baigiamojo magistro

### **darbo tikslas:**

Atlikus emocijų klasifikavimo metodų tyrimą, sukurti kuo efektyvesnį automatinį emocijų klasifikatorių, skirtą lietuviškų internetinių tekstų emocijų analizei.

Darbo tikslo įgyvendinimui keliami šie **uždaviniai**:

1. Išanalizuoti egzistuojančius sentimentų analizės metodus;
2. Atlikti mokslinių darbų sentimentinės analizės tematika apžvalgą;
3. Sugeneruoti emocijų klasių žodyną lietuvių kalbai;
4. Sudaryti tekstyną žodyno generavimui ir testavimui;
5. Įgyvendinti teksto emocijų klasifikavimo uždavinį semantiniu ir statistiniu metodais;
6. Pateikti teksto emocijų klasifikavimo sistemos kūrimo metodiką bei gaires tolimesniam šios srities tyrimų vystymui.

**Metodai.** Iškeltam tikslui ir uždaviniams įgyvendinti atlikta sentimentinės analizės metodų apžvalga, sudaryta emocijų vertinimo teorijos taksonomija. Teksto emocijų klasifikavimo algoritmas sumodeliuotas remiantis semantiniu ir statistiniu metodais. Sukurti algoritmai suprogramuoti Java kalba nenaudojant jokių trečiųjų šalių bibliotekų.

Statistinis metodas, paremtas prižiūrimuoju mašininu mokymu, įgyvendintas remiantis tekstynu su sužymėtais emocijų klasėmis dokumentais. Siekiant teksto emocijų klasifikatoriaus universalumo, sugeneruoti iš skirtingo žanro duomenų sudaryti naujienų ir komentarų tekstynai. Komentarų tekstynas skirtas apmokymui, o naujienų tekstynas – automatizuotu būdu sugeneruoto žodyno testavimui. Metodas paremtas chi-kvadratų ir regresijos algoritmais. Chi-kvadratų testas taikomas reikšmingų klasifikavimui požymių atrinkimui, tiesinės regresijos algoritmas – nustatytų reikšmingų požymių

svorių skaičiavimui. Sviurių skaičiavimo funkcijos optimizavimui pritaikytas gradientinio nusileidimo algoritmas.

Semantinio metodo įgyvendinimui sudarytas emocijas indikuojančių žodžių žodynas. Semantinis metodas paremtas emocijų žodyno ir lingvistinių taisyklių modelio taikymu. Žodynas sudarytas pagal emocijų vertinimo teorijomis ir konceptualiomis metaforomis pagrįstą taksonomiją. Lingvistinių taisyklių, įtakojančių teksto semantiką, modelis apima emocijos paneigimo, prasmės intensyvumo ir junglumo pokyčio bei ironijos nustatymo mechanizmus.

Atliktų eksperimentų metu buvo vertinama sistemos teksto emocijų klasifikavimo kokybė pagal pasirinktų metrikų rinkinį. Automatiniu būdu priskirta teksto emocijos klasė buvo lyginama su anotuotojų sužymėtomis klasėmis. Taip pat buvo lyginami sistemos apskaičiuoti teksto emocijos skaitiniai įverčiai su anotuotojų vertinimais. Pasirinktų metrikų rinkinys susideda iš keturių metrikų: tikslumas (angl. precision), klaidų kiekis (angl. error rate), atpažintų arba klasifikuotų objektų kiekis (angl. recall), F-įvertis (angl. F-score).

**Problematika.** Emocijų aptikimas yra sudėtingas uždavinys dėl emocijos išraiškos subtilybių, paralingvistinės informacijos, palengvinančios emocijos klasės nustatymą, trūkumo, todėl į emocijų klasifikavimo modelį integruoti ironijos, emocijos paneigimo ir sustiprinimo algoritmai.

Kita problema – automatizuotų teksto analizės sistemų rezultatų neobjektyvumas. Mašininio mokymo algoritmai itin derinami prie apmokymui skirtų duomenų, todėl sistema gerai identifikuoja sentimentus, kai testavimui skirti duomenys panašūs į apmokymui skirtus duomenis, tačiau tuo pačiu algoritmu gauti rezultatai ženkliai suprastėja, kai apmokymo ir testavimo duomenys skirtingi. Siekiant rezultatų objektyvumo apmokymui ir testavimui buvo naudojami skirtingos stilistikos duomenys.

Kuriant emocijų analizės sistemą problema yra jos pritaikymas lietuvių kalbai, nes tai specifinė kalba, pasižyminti sudėtinga gramatika, žodžių gausa.

Dar didesnė problema yra lietuvių kalbos vartojimas internete. Daug lietuvių yra atsisakę lietuviškų rašmenų ir nepaiso bendrinių lietuvių kalbos normų, bendraudami internete. Dalis komentarių autorių nenaudoja lietuviškų rašmenų. Nemaža dalis komentarių pasižymi rašybos klaidomis, todėl pirminiame duomenų paruošimo etape šias klaidas reikia ištaisyti, kad analizuojamas tekstas būtų parašytas bendrine kalba, kurią suprastų sistema. Kuriant emocijų analizės sistemą, būtina nustatyti lietuvių kalbos žodžių bendrines kalbos formas pasitelkiant trečiųjų šalių lietuvių kalbos morfologinės analizės ir žodžių suvedimo į pagrindinę formą įrankius.

**Struktūra.** Darbas yra tęstinis, jame tęsiamas mokslo tiriamojo darbo pradėtas eksperimentas. I skyriaus 1.1, 1.2 ir 1.3 poskyriuose pateikiama informacija yra aptarta mokslo tiriamajame projekte, tačiau šiame darbe atnaujinta. Tačiau eksperimentinė dalis, išskyrus naujienų tekstyno generavimą, atnaujinta, sumodeliuoti ir įgyvendinti nauji algoritmai Java programavimo kalba.

Darbą sudaro įvadas, trys pagrindiniai skyriai, rezultatų apibendrinimas, naudotos literatūros sąrašas ir ateities tyrimų gairės.

Įvadiniamame skyriuje aptariama tiriamoji problema, darbo aktualumas, aprašomas tyrimų objektas, formuluojamas pagrindinis darbo tikslas bei uždaviniai, aprašoma tyrimų metodika, darbo mokslinis naujumas, pasiektų rezultatų praktinė reikšmė. Įvado pabaigoje pristatoma magistrinio darbo struktūra.

Pirmajame skyriuje apibrėžtas ir detalizuotas sprendžiamas uždavinys, pateikta analitinė kitų autorių darbų apžvalga. Pasirinkti ir išanalizuoti keli populiarūs sentimentų analizės metodai, kurie eksperimentinėje darbo dalyje lyginti su autorės pasiūlytaisiais.

Antrajame skyriuje sudaryti semantinis ir statistinis internetinių dokumentų emocinio klasifikavimo modeliai, suformuluota metodika lietuviškų tekstų emocijų analizei, kuria remiantis darbo eigoje sukurtas įrankių rinkinys Java kalba.

Trečiajame skyriuje pateikti pagrindiniai eksperimentinio tyrimo rezultatai bei rekomendacijos, kurios leistų patobulinti sukurtą emocijų nustatymo tekste prototipą.



## 1. Analitinė dalis

### 1.1. Sentimentinės analizės taikymai teksto gavybos uždaviniuose

Sentimentinė analizė - tai automatizuotu būdu išgaunama nuomonių, vertinimų, sprendimų, poliariškumo faktų ir kitų subjektyvių teksto išraiškų visuma. Tai atliekama kompiuterinės lingvistikos ir mašininio mokymo metodų sintezės būdu.

Sentimentinė analizė priklauso teksto analizės užduotims.

Teksto analizė - tai įvairios (ne tik sentimentinės) informacijos nustatymas iš nestruktūrizuotų duomenų, t.y. teksto. Tam gali būti naudojama Python kalba parašyta pakankamai išsami NLTK biblioteka.

Sentimentinė analizė naudojama ne tik emocijų aptikimo, bet ir kitoms teksto analizės užduotims atlikti. Pastaraisiais dešimtmečiais sentimentinės analizės tyrimai atliekami įvairiais aspektais: subjektyvumo, sentimentų, emocijų. Nėra vieningo sentimentinės analizės klasifikavimo. Tačiau visų sentimentinės analizės krypčių ištakos – afektinė kompiuterija, todėl šiame darbe remiamasi Picard ontologija ir terminu „afektas“ apibūdinama emocijos, sentimentai, asmenybės, nuotaikos ir požiūriai.

Scherer [11] apibrėžė šias afektines būsenas remdamasis pažinimo ir laiko faktoriais:

Emocija – tai santykinai trumpa, bet stipri atsako į dirgiklį reakcija, sukelta staiga pakitusių svarbių subjektui gyvenimo aplinkybių.

Nuotaika – sklaidi afektinė būsena, labiausiai pasižyminti subjektyvumo pokyčiu. Tai žemo intensyvumo, bet santykinai ilgos trukmės procesas, kylantis be aiškios priežasties (paniuręs, džiaugsmingas, suirzęs).

Asmeninė pozicija – pasikeitusi afektinė pozicija kito asmens atžvilgiu (draugiškumas, šaltumas).

Požiūris – ilgalaikiai emocinio pobūdžio įsitikinimai ir vertybės (mėgimas, neapykanta, vertinimas).

Asmenybės bruožai – emocinio turinio pastovūs asmenybės ir elgsenos bruožai (nervingumas, žiaurumas, nepastovumas).

Šios afektinės klasės glaudžiai susijusios, todėl emociją galima apibrėžti afekto terminu.

Iš kitos pusės, sentimentai, emocijos, nuomonės taip pat glaudžiai susiję. Liu [21] nuomone, sentimentų, subjektyvumo ir emocijos konceptai nėra tapatūs. Sentimentai gali būti išreikšti objektyviais teiginiais, t.y. tiesiog faktais (pvz. „gaisras sunaikino mano darbą“), kuriuose nėra jokio subjektyvaus požiūrio. Kita vertus, emocijos gali būti reiškiamos tik per subjektyvumo prizmę, bet jos nebūtinai atspindi nuomonę (pvz. „Buvau nustebęs, matydamas kas vyksta“). Be to, galima išreikšti sentimentą ir nuomonę, nelydimą jokios konkrečios emocijos (pvz. „Įmonės vadovai yra patyrę darbuotojai“). Visus šiuos ypatumus būtina kruopščiai apsvarstyti siekiant nustatyti tikslias sentimentinės analizės ribas. Šiame tyrime sentimentai apibrėžiami kaip bet kokia asmeninė subjektyvi nuomonė, kuri gali būti reiškiamas ar neigiamas emocijas, vertinimus ir pozicijas. Žemiau lentelėje pateikiami sentimentų raiškos pavyzdžiai.

**1 lentelė.** Sentimentų raiškos aspektų pavyzdys

	Teigiami sentimentai	Neigiami sentimentai
Emocija	Aš laiminga	Jis liūdnas
Vertinimas	Puiki mintis!	Nekokia mintis!
Pozicija	Ji už įstatymą	Jis prieš įstatymą

1 lentelėje pateikti sentimentų raiškos per emocijas, vertinimus ir pozicijas pavyzdžiai rodo glaudžią įvairių subjektyvumo aspektų ir afektinių klasių sąsają.

## **1.2. Emocijų analizė – sentimentinės analizės kryptis**

### **1.2.1. Emocijų analizės raida**

Viena svarbiausių afektinių klasių – emocijos. Emocijų nustatymas atlieka reikšmingą vaidmenį daugumoje teksto analizės uždavinių ir yra taikomas įvairiose srityse.

1. Visuomenės sveikata: Itin svarbi emocijų atpažinimo įtaka medicinoje depresijos [12], polinkio į savižudybę [17], kibernetinio smurto atpažinimui [13], bendruomenės sveikatos būklės, arba geros savijautos nustatymui [15]. Taip pat atliekami eksperimentiniai robotų-pagalbininkų projektavimo darbai, siekiant sukurti robotus, galinčius nustatyti pagyvenusius, neįgalių ar sergančių žmonių emocijas ir pagal poreikį suteikti jiems fizinę terapiją paslaugas.

2. Politika: Itin didžiulis domėjimasis visuomenės sentimentais, susijusiais su politika, ypač pranašaujant rinkimų rezultatus [11].

3. Vadyba: Elektroninių dienoraščių, "Twitter", "Facebook" komentarų analizė plačiai taikoma prekių ženklų rinkodaroje, klientų nuomonės, vartojimo tendencijų tyrimuose. Automatizuotas emocijų nustatymas klientų atsiliepimuose padeda verslui palaikyti grįžtamąjį ryšį su klientais, išsiaiškinant prekių ar paslaugų trūkumus bei privalumus [23].

4. Švietimas: Studentų vertinimo sistemose pagal emocijas nustatyti atsakymų teisingumą ir emocinę atsakinėjančiojo būseną. Emocijų nustatymas kompiuterizuoto mokymo sistemose padeda išsiaiškinti kokias emocijas mokomoji medžiaga sukelia besimokančiajam [22].

5. Asmenybės bruožų nustatymas: Pagal emocijų reiškimo pobūdį nustatyti asmenybės bruožai tokie kaip ekstravertiškumas ir narcisizmas [11].

6. Literatūros analizė: Vis didėja susidomėjimas natūralios kalbos priemonių panaudojimu didelės apimties literatūrinių tekstų analizei [19].

7. Neuromarketingas: Verslui vis aktualesnė tampa rinkodaros kryptis, tirianti žmogaus neurologines reakcijas į reklamas ir produktus – tyrimai šiandien įgalina išmatuoti bazines vartotojo emocijas bei jas lydinčias fiziologines reakcijas.

Teksto emocinio fono nustatymas buvo atliekamas daugumoje įvairių dokumentų: pasakose ([8], [27]); elektroniniuose dienoraščiuose ([29], [19]), romanuose [9], pokalbių pranešimuose [18] ir socialinės žiniasklaidos dokumentuose [12]. Lyginamoji emocijų žodžių pasiskirstymo elektroniniuose dienoraščiuose ir asmeninio turinio dokumentuose (meilės, savižudybės laiškuose) analizė atskleidė, kad emocijos, tokios kaip pasibjaurėjimas, vienodai išreiškiamos socialinio tinklo „Twitter“ žinutėse ir meilės laiškuose [10].

Naujienu portalu emociju analizės pradžia – SemEval “Affective Text“ konkursas 2007 metais. Naujienu žinučių emociju klasifikavimą žodžių krepšelio metodu pirmieji įgyvendino P. Katz, M. Singleton ir R. Wicentowski [8], C. Strapparava ir R. Michalcea pritaikė asociacijų metrikos (PMI) bei sintaksės taisyklių ir žodyno metodus.

### **1.2.2. Emocijų aptikimo tekste problematika**

Emocijos dominuoja visose gyvenimo srityse. Jos įtakoja sprendimų priėmimą, santykius, apsprendžia mūsų elgesį. Žmonės reiškia emocijas įvairiais tiek verbaliniais, tiek neverbaliniais būdais: kalbos išraiškomis, balso intonacijomis, kūno kalbos signalais. Nors atlikta nemažai emocijų aptikimo iš balso, veido išraiškų ir fiziologinių signalų tyrimų, tačiau emocijų aptikimo iš rašytinio teksto sritis yra pastangų ir įžvalgų reikalaujantis uždavinys. Išsakant emocijas kalboje vyrauja ekspresinė (lot. *expressio* – išreiškimas, išraiška), arba emocinė, funkcija, atspindinti autoriaus santykį su teksto turiniu ir adresatu. Atpažįstant emociją tekste analizuojama kalbėtojo vidinė būseną, vertinimai, o ne teksto turinys.

Socialiniuose tinkluose ir diskusijų forumuose emocijų nustatymo užduotis palengvina sutartiniai emocijų reiškimo jaustukai (angl. *emoticons*), grotelių (angl. *hashtags*) ir grafiniai (angl. *emoji*) simboliai, tačiau emocijų nustatymo uždavinį apsunkina necenzūrinės išraiškos, barbarizmai, piktogramos, nelietuviški rašmenys, gramatinės klaidos. Naujienu portalu žinučių privalumas – bendrinės kalbos naudojimas, tačiau trūkumas – kalbos ekspresyvumo stoka.

Nepriklausomai nuo diskurso emocijų analizės tyrėjai susiduria ir su kitais iššūkiais, susijusiais su subjektyvumo modeliavimo ir emocinio turinio sudėtingumo problematika:

1. Sentimento išraiškos subtilybės:
  - Perkeltinės kalbos naudojimas. Sentimentui reikšti naudojamos įvairios retorinės figūros: ironija, sarkazmas, hiperbolė, metaforos, similės;
  - Nuomonė išreiškiama neutraliais žodžiais.
2. Neiginių ir tikimybinių išraiškų naudojimas gali pakeisti sakinio sentimentą, pvz. geras, negeras ir turėtų būti geras turi visiškai skirtingas kanotacines reikšmes.
3. Srities/konteksto prigimtis. Tie patys žodžiai arba frazės gali reikšti skirtingus dalykus skirtingame kontekste.
4. Specifiniai socialiniuose tinkluose naudojami sintaksiniai dariniai. Socialinės žiniasklaidos tekstuose gausu terminų, kurių nėra jokiuose žodynuose: neteisingos rašybos naujadarai (happee), žodžiai su grotelių simboliu (#loveumom), jaustukai, grafiniai simboliai (angl. *emoji*), sutrumpinimai (4U). Paprastai šie terminai išreiškia emocijas.
5. Paralingvistinės informacijos (balso intonacijos, kūno kalbos) stoka.
6. Anotuotų duomenų trūkumas. Emocijų yra įvairių, bet prieinama tik 6-8 emocijomis anotuoti duomenų rinkiniai.
7. Automatizuotų teksto analizės sistemų rezultatų neobjektyvumas. Skirtingų tyrėjų sukurtais algoritmais grindžiamos automatizuotos teksto analizės sistemos paprastai vertinamos, naudojant skirtingus testavimo duomenis ir nustatymus. Todėl atskiruose straipsniuose pateiktų rezultatų, negalima objektyviai palyginti. Taigi negalima nustatyti kuris metodas geriausiai tinka praktiniam

pritaikymui. Be to, mašininio mokymo algoritmai itin derinami prie apmokymui skirtų duomenų, todėl sistema gerai identifikuoja sentimentus, kai testavimui skirti duomenys panašūs į apmokymui skirtus duomenis, tačiau tuo pačiu algoritmu gauti rezultatai ženkliai suprastėja, kai apmokymo ir testavimo duomenys skirtingi.

8. Subjektyvūs ir tarpkultūriniai skirtumai.

### 1.2.3. Emocijų vertinimo taksonomija

Emocijų analizė – tai kompiuterizuotas natūralios kalbos išraiškų tyrimas, siekiant šioms išraiškoms priskirti įvairias emocijas.

Kompiuterinėje emocijų analizėje naudojamos emocijų kategorijos, grindžiamos psichologijos ir pažinimo teorijos mokslais, o natūralios kalbos apdorojimo ir teksto analizės metodai apdoroja tekstą ir pagal nustatytas emocijų kategorijas įvertina teksto emocinį foną.

Svarų indėlį į emocijų tyrimą psichologijoje įnešė Ortony A., Clore G.L. ir Collins A. [19]. Jų teorija, vadinama OCC emocijų modeliu, turėjo įtakos kompiuterinei lingvistikai [15]. Akronimas OCC sudarytas iš autorių pavardžių pirmųjų raidžių.

Visi emocijų klasifikavimo modeliai, padedantys išvelgti ir vertinti žmonių emocines išraiškas, grindžiami psichofiziologine emocijų kilme, tačiau kiekvienas jų pasižymi savita emocijų vertinimo sistema. Emocijų analizės uždaviniai paremti trimis plačiausiai emocijų analizėje naudojamomis emocijų vertinimo teorijomis.

- Diskrečioji emocijų klasifikavimo teorija yra paremta fiziologine atsako į dirgiklį kilme: išskiriamos diskrečios emocijos, sužadinančios skirtingus fiziologinius procesus. Populiariausias diskrečios, kitaip dar vadinamos universalios arba atominės, emocijų vertinimo teorijos atstovas – Ekman, suskirstęs emocijas pagal fiziologinės reakcijos pobūdį į šešias kategorijas: laimė, pyktis, nuostaba, liūdesys, baimė ir pasibjaurėjimas. Kita diskrečiosios krypties Plutchik teorija klasifikuoja emocijas į aštuonių kategorijų keturias priešingų emocijų poras: laimė-liūdesys, pyktis-baimė, pasitikėjimas-pasibjaurėjimas, numatymas-netikėtumas;
- Dimensinė emocijų klasifikavimo teorija vertina emocijas dviejų arba trijų dimensijų erdvėje. Pagrindiniai vertinimo kriterijai: junglumas, susijaudinimo laipsnis ir dominavimas [16]. Junglumu (angl. *valence*) apibūdinamas dirgiklio sukeliamas emocijos vertinimo poliariškumas. Susijaudinimo laipsniu (angl. *arousal*) žymimas emocijos intensyvumas. Dominavimu (angl. *dominance*) matuojamas sužadintos emocijos valdymo laipsnis;
- Vertinamoji arba, kitaip prototipo, teorija paremta subjektyvumo idėja, aiškinančia skirtingas to paties prototipo šeimos emocijas kaip atsako į dirgiklį priklausomybę nuo reaguojančio organizmo pajėgumo, pvz. esant pavojaus dirgikliui išsekęs atsitrauks nuo pavojaus, stiprus – puls.

Neseniai buvo susidomėta bazinio emocijų karkaso papildymu kiek sudėtingesnėmis emocijomis tokiomis kaip mandagumas, grubumas, apgaulė, depresija, gyvybingumas ir sumaištis [22]. Bendra visiems emocijų karkasams savybė – emociškai stiprių išraiškų parinkimas ir emocijos stiprumo nustatymas [23]. Atsižvelgiant į šias teorijas, sudaromi emocijas indikuojančių žodžių sąrašai, palengvinantys emocijos klasės tekste nustatymą. Tokie sąrašai vadinami emocijų žodynais, arba leksikonais.

#### 1.2.4. Emocijų raiškos internetinėje erdvėje būdai

Verbalinės emocijų raiškos priemonės – žodžiai – realiame gyvenime nėra svarbiausias emocijos indikatorius. Kai kurių tyrėjų nuomone [8], žodžiai sudaro tik 10 proc. emocijos raiškos. Daug svarbesni emocijos žymikliai – veido išraiškos, kūno judesiai, balso tembras, todėl rašytinio teksto emocijų klasifikavimo uždavinys, remiantis tik verbaliniais emocijos žymikliais, nėra lengvas. Emocijų nustatymo palengvinimui pasitelkiamos įvairios neverbalinės priemonės: grotelių (angl. *hashtag*), grafiniai (angl. *emoji*), sutartiniai emocijų reiškimo (angl. *emoticons*) simboliai ir nuostatų skalės. Šie emocijų žymikliai indikuoja dokumento emociją, kurios klasei dokumentas gali būti priskiriamas automatiškai be rankinio anotavimo, todėl vis dažniau taikomi teksto emocijų klasifikavimui. Kinijos naujienų portaluose skaitytojams suteikiama galimybė įvertinti straipsnio sužadintą emociją, pažymint vieną iš aštuonių emocijų reiškimo simbolių.

Nuostatų skalės, leidžiančios pažymėti vertinimą, labiau taikomos sentimento identifikavimui. Nuostatų skalės nustato požymio raiškumo laipsnį, fiksuoja nuostatas į procesus ar reiškinius. Internetinių portalų komentaruose leidžiamas tik binarinis vertinimas „patinka“ arba „nepatinka“, tačiau atsiliepimų apie prekes ir paslaugas puslapiuose galima vertinti balais arba suteikiant tam tikrą kiekį žvaigždučių. Vertinimo skalėmis, grafine simbolika išreikštą emociją analizuoti yra kur kas lengviau negu emociją, išreiškiamą tik žodžiais.

Tačiau Lietuvos naujienų portalų pranešimuose, paklūstančiuose bendrinės kalbos reikalavimams, sutartiniai emocijų raiškos, grotelių ar grafiniai simboliai neleidžiami. Bendrine kalba parašytame tekste emocijoms reikšti naudojamos sintaksinės priemonės, pavyzdžiui pyktis atpažįstamas iš daugtaškio, reiškiančio nutylėjimą arba šauktuko, reiškiančio retorinį sušukimą, nuostaba – iš klaustukų ar klaustuko ir šauktuko derinio. Pasitelkiant sintaksines priemones modeliuojamos retorinės figūros, didinančios kalbos ekspresyvumą, kurioms būdingas žodžių vartojimas perkeltine prasme, aiškinant vieną reiškinių kitu.

Viena svarbiausių retorinių struktūrų – metafora – atlieka svarbų vaidmenį emocijos žodyno generavimui, o tam tikromis metaforos atmainomis (ironija, hiperbolė ir t.t.) paremtas emocijų nustatymo modeliavimas, todėl šiame skyriuje pateikiama glausta meninių raiškos priemonių, vadinamų retorinėmis figūromis, naudotų emocijos analizės prototipe, samprata.

Pagal reikšmės perkėlimo būdą skiriamos dvi retorinių figūrų atmainos:

- 1) metafora – kai reikšmės perkėlimo pagrindą sudaro vaizdo panašumas,
- 2) metonimija – kai reikšmės perkėlimo pagrindą sudaro loginis sąvokų ryšys. Dažnai šios dvi atmainos nėra grynos, o persipynusios tarpusavyje, todėl tapatinamos. Toliau metafora vadinama bet kuri iš šių atmainų. Pavartoti metaforą – tai reiškia pastebėti, kas panašu tarp dviejų nieko bendra neturinčių reiškinių, įmanomų iš tolimiausių ir skirtingiausių pasaulių. Šia savybe paremta konceptualių metaforų idėja, kuria remtasi sudarant emocijų žodyną.

Konceptualios metaforos – vienas esminių mąstymo ir pažinimo būdų, formuojančių mąstymą, suvokimą, veiksmą. Konceptualių metaforų teorija paremta kognityvinės lingvistikos mokslu, tiriančiu kalbos-minties-kūno ryšį. [Gibbs]

Modeliuojant emocijų analizės prototipą, atsižvelgta į šias metaforos atmainas:

- ironija – reikšmės perkėlimas priešingą reikšmę turinčiam žodžiui. Jai skirtas ypatingas dėmesys, sukurta autonomiška (atskira) ironijos nustatymo tekste posistemė;
- hiperbolė – sąmoningas perdėjimas, ką nors apibūdinant – tiesiogiai sietina su ironijos atpažinimu rašytiniame tekste;
- oksimoronas – antoniminių žodžių gretinimas;
- pakartojimas – žodžių, junginių ar sakinių kartojimas, siekiant sustiprinti kalbos raiškumą, pabrėžti reiškinio svarbą, veiksmo ar būsenos trukmę, intensyvumą;
- nutylėjimas – stilistinę paskirtį turintis staigus kalbos nutraukimas, paliekant pačiam skaitytojui suvokti, kas nepasakyta;
- retorinis sušukimas – šaukiamasis sakiny, kuris ypač emociškai iššreiškia mintį;

Retorinės figūros sumodeliuotos pasitelkiant verbalines – emocinius žodžius ir frazes bei neverbalines – emocijų reiškimo simbolius, skyrybos ženklus – priemones.

**2 lentelė.** Emocijų raiškai naudotos stilistinės priemonės

Retorinė figūra	Verbalinės priemonės	Neverbalinės priemonės
Ironija	frazės, žodžiai	hiperbolė, oksimoronas, pakartojimas
Hiperbolė	žodžiai, jaustukai	emocijų reiškimo simboliai
Oksimoronas	žodžiai	
Nutylėjimas		Skyryba (daugtaškis)
Retorinis sušukimas		Skyryba (šauktukas)

2 lentelėje pateiktų emocijų raiškai naudojamų retorinių ir sintaksinių priemonių pavyzdžiai patvirtina, kad emocijų analizės procesas paremtas psichofiziologijos, lingvistikos ir kompiuterių mokslo metodų sinteze.

### 1.3. Sentimentų žodynai

#### 1.3.1. Sentimentų žodynų apžvalga

Emocijų analizės, kaip ir bet kurios kitos sentimentinės analizės, užduočiai atlikti būtinas sentimentų žodynas. Minimali tokio žodyno funkcija - automatizuoto klasifikavimo įvertinimas, tačiau, taikant mašininio mokymo metodus, žodynai gali būti naudojami ir kaip apmokymui skirti duomenys [1]. Kadangi sentimentas atspindi emociją, o emocija atspindi požiūrį, kuris tapatinamas su afektine būsena, tai darbe sutinkami terminai emocijų, sentimentų arba afektinis žodynas naudojami lygiagrečiai emociją indikuojančių žodžių sąrašui apibūdinti. Dažniausiai tokiuose žodynuose žodžiai sužymimi vienos dimensijos (sentimento, junglumo, asociacijos) binariniais įverčiais, žyminčiais dimensija vertinamos savybės faktą.

Pirmasis sentimentų žodynas, sudarytas iš 3600 pagal poliariškumą sužymėtų žodžių – General Inquirer [Stone], paremtas kognityvine žodžio prasmės vertinimo teorija.

MPQA subjektyvumo žodynas [2] yra vienas geriausiai vertinamų žodynų anglų kalba. Jis sudarytas iš 2718 teigiamų ir 4912 neigiamų žodžių, surinktų iš įvairių šaltinių: General Inquirer sąrašų, Hatzivassiloglou ir McKeown asociacijų pagrindu sugeneruoto žodyno ir rankiniu būdu sudaryto subjektyvumą žyminčių žodžių ir frazių sąrašo [6]

LIWC, Linguistic Inquiry and Word Count, yra 73 žodžių sąrašų, sudarytų iš 2300 žodžių, rinkinys [17], skirtas socialinės psichologijos užduotims.

Poliariškumo leksikonas, sudarytas iš 2006 teigiamų ir 4783 neigiamų žodžių – pirmasis žodynas, paremtas žodžių atranka iš socialinio turinio teksto – klientų atsiliepimų [23]

Visi šie žodynai yra anglų kalba. Keletas sentimentų žodynų yra prieinami kitomis kalbomis, pvz., MLSA [6] yra pirmasis viešai publikuotas šaltinis vokiečių kalba. Chen ir Skiena [5] nustatė 12 viešai publikuojamų žodynų penkiomis kalbomis; tačiau tarp jų nėra nei vieno lietuvių kalba. Nors 2013 m. Radovan Garabik ir Indre Pileckyte paskelbė sentimentų žodyną slovakų-lietuvių kalbomis, tačiau jame nėra emocinio žymėjimo[30].

### **1.3. 2. Sentimentų žodynų generavimo metodai**

Emocijų, kaip ir kitų sentimentinės analizės aspektų, žodynų generavimui taikomi šie metodai:

1. Rankinis būdas. Visi populiariausi žodynai sudaryti rankiniu būdu. Tai – The General Inquirer, Hu ir Liu poliariškumo žodynas, MPQA Subjectivity Lexicon ir pripažinimo sulaukę emocijų žodynai – Strapparava ir Valitutti Wordnet Affect [25] ir Saif M. Mohammad NRC[17]. Pastaruoju metu rankinis būdas modernizuojams, taikant automatizuoto masinio anotavimo (angl. *crowdsourcing*) metodą, pasitelkiant saityno paslaugą. Naudojant masinio anotavimo metodą, emocijos priskyrimo žodžiams užduotis paskirstoma daugumai anototojų, ir vėliau apibendrinant rezultatus, žodžiui priskiriama emocijų klasės, o neretai ir keletos klasių žyma, atitinkanti daugumos nuomonę. Tokiu būdu sudarytas NRC Word-Emotion Association Lexicon, kuriame 8 emocijų klasėmis pagal Plutchik taksonomiją ir poliariniais sentimentais sužymėta apie 14,000 žodžių[21]. Žymėjimas buvo atliekamas naudojant žiniatinklio paslaugą Mechanical Turk.

Kitas žodynas [29], sudarytas naudojant masinio anotavimo metodą, paremtas trijų dimensijų – junglumo, sujaudinimo, dominavimo – įverčių priskyrimu analizuojamiems žodžiams. Sužymėta 14,000 žodžių kiekvienai nustatytų dimensijų priskiriant įverčius pagal 1-9 nuostatų skalę.

2. Kadangi rankiniu būdu sudaryti žodyną reikalauja daug laiko ir žmogiškųjų resursų, tampa įprasta žodynus generuoti automatizuotai taikant žodyno metodą. Tokių žodynų privalumas – didesnė apimtis. Pavyzdžiui, SentiWordNet [20] susideda iš 38,182 afektinių žodžių, o Maryland žodynas [19] iš 76,775 žodžių ir frazių su poliariškumo žymomis. Nors jų apimtis didesnė lyginant su rankiniu būdu sudarytais žodynais, tačiau dėl neišsamiai pateiktos informacijos automatizuotu būdu sudaryti žodynai pasižymi žemesne kokybe.

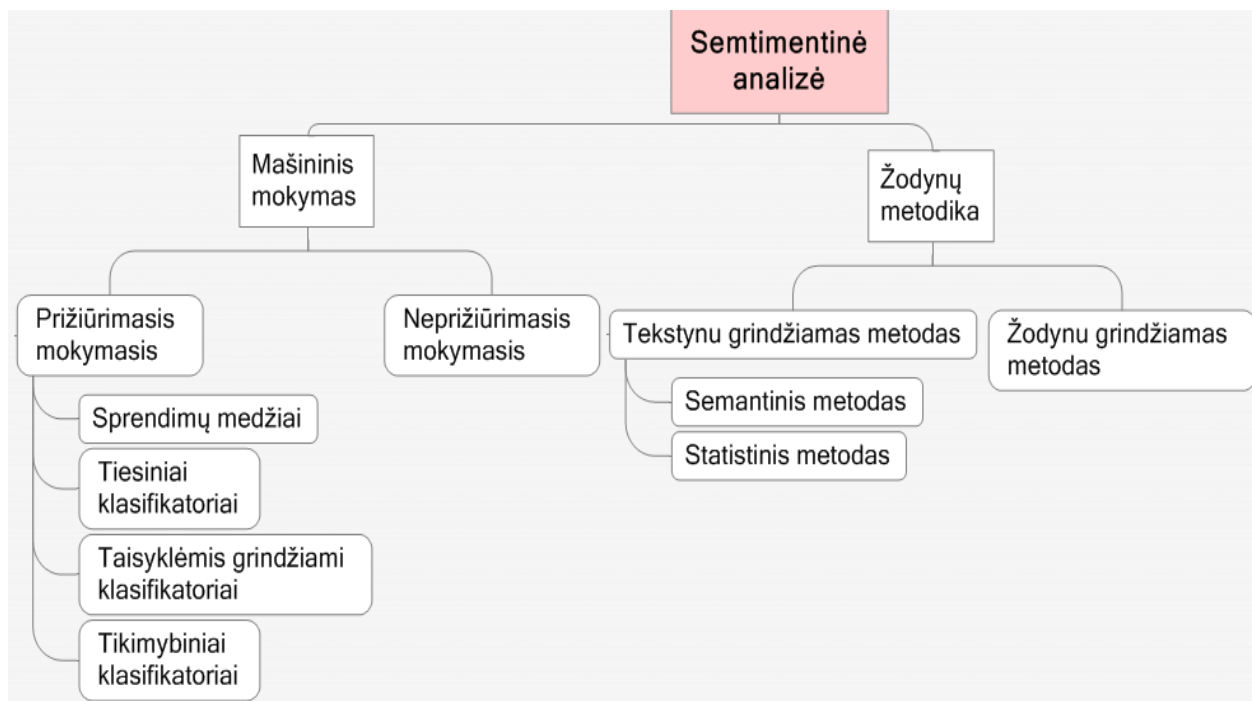
Žodynu grindžiamo metodo eiga – keletui sentimentinių žodžių internetiniame žodyne ieškomi sinonimai ir antonimai. Paprastai remiamasi WordNet žodynu dėl šio žodyno išsamaus turinio ir taikymo įvairiose srityse. Kamps naudojo WordNet atstumu grindžiamą metodą būdvardžio sentimentinės orientacijos nustatymui [16], Williams ir Anand rėmėsi sentimentų stiprumo savybe [14]. Alternatyva WordNet žodynui – tezaurai, pvz., naudodamiesi tezauru Mohammad, Dunne ir Dorr automatizuotai sugeneravo sentimentų žodyną, susidedantį iš daugiau nei 60,000 žodžių [26]. Tačiau neformali ir dinamiška internetinės erdvės prigimtis apsunkina šių žodynų pritaikymą emocijų analizei.

Alternatyva – emocinių savybių išgavimas iš dinaminio, t.y. apmokyto sentimentų žodyno. Šiuo atveju emocinio poliariškumo tarp žodžio ir emocijomis sužymėto turinio modeliavimui naudojama asociacijų metrika (angl. *Pointwise Mutual Information – PMI*) [28].

3. Hibridinis, arba tekstynu grindžiamas metodas: turint pradinį bendrojo pobūdžio sentimentinių žodžių rinkinį ir taikant mašininio mokymo metodus, ieškoma kitų sentimentinių žodžių srities kontekste arba adaptuojant bendrosios paskirties sentimentų žodyną tiriamai dalykinei sričiai. Hibridiniu metodu sugeneruoti žodynai, su statinių emocinių žodžių pagrindu išgautais dviejų-trijų žodžių junginiais pasižymi tikslesniu teksto emocinio fono nustatymu negu žodynai tik su žodžiais [26].

Populiarėjant socialiniams tinklams, prigijo nuotolinio mokymo metodas, kuriuo rankinį tekstyno žymėjimą pakeičia grotelių ir emocijų reiškimo simboliams priskiriamos klasių žymos. Dažniausiai dėl programavimo sąsajos patogumo ir teksto lakoniškumo analizuojamos Twitter socialinio tinklo žinutės. Kiritchenko, Zhu, Mohammad iš Twitter žinučių automatizuotai sugeneravo emocijų žodyną, kurį integravę į teksto emocijų klasifikavimo sistemą, pasiekė 87 proc. tikslumą ir laimėjo SenEval-2013 konkursą. [15] Go, Bhayani ir Huang naudojo Twitter žinučių emocijų reiškimo simbolius kaip poliariškumo klasių žymas, taikydami prižiūravimo mokymo metodą. Mohammad sukūrė klasifikatorių Twitter žinučių emocijų nustatymui naudodamas grotelių simbolius su klasių pavadinimais (pvz., #anger, #surprise).

### 1.3. Sentimentinės analizės metodų apžvalga



1 pav. Klasifikavimo metodų schema



Kadangi emocijos aptikimas tekste yra laikomas vienu iš sentimentinės analizės aspektų, todėl emociinei teksto analizei taikomi bendri sentimentų klasifikavimo metodai. Išskiriamos trys pagrindinės sentimentų klasifikavimo metodikos:

- Mašininio mokymo metodika
- Žodynu grindžiama metodika
- Hibridinė metodika[69].

Mašininis mokymas ( angl. Machine Learning – ML) – tai procesas, kurio metu informacija paverčiama žiniomis, užuot iš anksto tas žinias įdiegus į sistemą. Sentimentinės teksto analizės atveju taikomi populiarūs mašininio mokymo algoritmai, apmokymui naudojantys lingvistinius požymius. Žodynu grindžiama metodika paremta sentimentinio žodyno naudojimu. Sentimentinis žodynas – tai iš anksto sudarytas žinomų sentimentinių terminų rinkinys. Hibridinė metodika apjungia mašininio mokymo ir žodyno metodus, sentimentų žodynui priskirdama pagrindinį vaidmenį teksto klasifikavimo uždaviniuose.

### **1.3.1. Klasifikavimui naudingų požymių atranka**

Visų metodikų pagrindas – klasifikavimo požymių atrinkimas ir interpretavimas. Informatyvūs klasifikavimui požymiai teksto analizės uždaviniuose gali būti aprašomi žodžių krepšeliu, modeliais arba ontologijomis.

Mašininio apmokymo atveju klasifikatoriai naudoja žodžius, frazes ar kitas savybes kaip klasifikavimo požymius. Taikant žodyno metodą, žodžiams ar lingvistinėms taisyklėms priskiriamas emocijos klasės, kuriai jie priklauso, įvertis. Žodžiai atrenkami iš sužymėto žodyno, taisyklės – modeliuojamos, kategorijos aprašomos ontologijomis – atrinktų žodžių pagal pasirinktą taksonomiją ir sumodeliuotų taisyklių sąrašais. Šie požymiai atrenkami taikant įvairius kompiuterinius metodus.

Viena teksto klasifikavimo problemų yra požymių erdvės daugiadimensiškumas, todėl kruopštus klasifikavimui naudingų požymių atrinkimas yra svarbus teksto klasifikavimo uždaviniuose. Požymių atrankos metodais siekiama iš požymių aibės išrinkti tokį požymių poaibį, su kuriuo būtų galima pasiekti maksimalų efektyvumą, tuo pačiu stengiantis neprarasti klasifikavimo tikslumo.

Emocijų klases indikuojančių reikšmingų požymių atrinkimui naudojami tiek rankinis būdas, tiek mašininio mokymo metodai.

Rankinis būdas grindžiamas taksonomijų arba anotavimo sistemų taikymu. Vienas iš anotavimo metodų – MaxDif, grindžiamas analizuojamų tarpusavio žodžių palyginimu [12]. Rankiniam žodžio emocijos klasės priskyrimui gali būti pasitelkiama Amazon saityno paslauga Mechanical Turk, įgalinanti masinį žodžių anotavimą.

Rankinis būdas reikalauja daug laiko, todėl paprasčiau reikšmingų požymių atrankai naudoti statistinius metodus. Nustatyta, kad efektyviausi požymių atrankos metodai – chi-kvadrato ir informacijos įgijimo [19].

- Paprasčiausias atrankos metodas – dokumentų dažnių metodas (Document frequency –DF). Taikant šį metodą požymiai, kurie pasitaiko rečiau nei nustatytą minimalų kartų kiekį, pašalinimi.
- Chi-kvadrato ( $\chi^2$ ) metodas yra grindžiamas statistine teorija ir vertina nepriklausomumą tarp požymio ir klasės [3].

- Informacijos įgijimo (angl. *Information Gain – IG*) metodu apskaičiuojama entropijos sumažėjimo tikimybė, skaičiuojant požymio aptikimo faktą ir susijusios klasės pasiskirstymo informaciją.

### 1.3.2. Mašininio mokymo metodų apžvalga

Mašininio mokymo metodika yra paremta mašininio mokymo algoritmų taikymu naudojant lingvistinius požymius. Teksto klasifikavimas, pagrįstas ML, skirstomas į prižiūravimo ir neprižiūravimo apmokymo metodus.

Prižiūravimo apmokymo atveju naudojami sužymėti pagal klases apmokymo duomenys, neprižiūravimo – nesužymėti duomenys.

Mašininio mokymo algoritmai atlieka teksto klasifikavimo uždavinį, naudodami teksto sintaksines ir/arba semantines ypatybes. Sekančiuose skyreliuose pateikiama automatizuoto apmokymo metodų apžvalga. Detalesnis eksperimente naudojamų algoritmų aprašymas pateikiamas antroje darbo dalyje, aprašančioje klasifikavimo eksperimentą.

#### 1.3.2.1. Prižiūrimasis mokymas

Teksto klasifikavimo užduoties apibrėžimas: naudojama apmokymui skirtas dokumentas  $D = \{X_1, X_2, \dots, X_n\}$ , kuriame kiekvienam įrašui  $X_1 \dots X_n$  priskirta atitinkamos klasės žymė. Klasė priskirta atsižvelgiant į tai klasei būdingų požymių buvimą įrašė. Tada apmokytas modelis naudojamas nesužymėto teksto klasifikavimui. Teksto klasifikavimo užduotis sudėtingesnė, kai apmokymui skirti duomenys sužymimi tik priskiriant klasę, tačiau šią užduotį palengvina priskirtos klasės ar kelių klasių tikimybinį įverčių žymėjimas. Yra keletas automatizuotam apmokymui skirtų klasifikatorių tipų.

*Tikimybiniai klasifikatoriai.* Tikimybiniai klasifikatoriai naudoja mišrius klasifikavimui skirtus modelius. Mišriu modeliu vadinamas modelis, sudarytas iš įvairių klasių komponentų, generuojančių atitinkamo emociją indikuojančio požymio priskyrimo atitinkamai klasei tikimybę.

Populiariausi tikimybiniai klasifikatoriai – Naivusis Bayeso, Bayeso tinklas, Maksimalioji Entropija.

Naivusis Bajeso klasifikatorius remiasi Bayeso tikimybių taisykle. Laikoma, kad visi duomenų požymiai yra nepriklausomi, ir kiekvienas iš požymių daro įtaką klasifikavimo rezultatui. Klasifikatorius skaičiuoja aposteriorines (angl. *posterior*) tikimybes kiekvienai klasei. Objektas priskiriamas tai klasei, kuri įgyja didžiausią aposteriorinę tikimybę [26].

Nustatant dokumento  $X$  kategoriją yra skaičiuojama klasės aposteriorinė tikimybė arba hipotezė  $h$ ,  $P(h|X)$ , kuri yra išreiškiama formule:

$$P(h|X) = \frac{P(X|h)P(h)}{P(X)}, \quad (1)$$

kur  $P(h) = \frac{|h|}{N}$  yra  $h$  apriorinė tikimybė ( $|h|$  ir  $N$  atitinkamai yra dokumentų skaičius klasėje  $h$  ir dokumentų skaičius visose klasėse, darant prielaidą, kad visos hipotezės yra vienodai tikėtinos),  $P(X|h)$  –  $X$  aposteriorinė tikimybė, kurią sąlygoja  $h$ , o  $P(X)$  –  $X$  apriorinė tikimybė, lygi konstantai. Siekiant sumažinti skaičiavimo kaštus, klasifikatorius daro naivią ir supaprastintą prielaidą, kad  $n$

atributų yra nepriklausomi vienas nuo kito. Sakykime, kad egzistuoja  $C$  klasių  $c_1, c_2, \dots, c_{|C|}$ , klasifikatorius daro prielaidą, kad nežinomas dokumentas  $x$  priklauso klasei, turinčiai didžiausią aposteriorinę (angl. *a posteriori*) tikimybę:

$$\arg \max_c \frac{P(c)P(\vec{x} | c)}{P(\vec{x})} = \arg \max_c P(c)P(\vec{x} | c) = \arg \max_c P(c) \prod_i P(x_i | c) \quad (2)$$

$\vec{x}$  – požymių vektorius,  $c$  -klasė [25].

*Tiesiniai klasifikatoriai.* Tarkime,  $\vec{X} = \{x_1, \dots, x_n\}$  yra normalizuotas dokumento žodžių dažnis,  $\vec{A} = \{a_1, \dots, a_n\}$  yra tiesinių koeficientų vektorius to paties dimensiškumo kaip ir požymių sritis, o  $b$  yra skaičius. Tada tiesinis prediktorius  $p = \vec{X} \cdot \vec{A} + b$  yra tiesinio klasifikatoriaus rezultatas. Prediktorius  $p$  yra skiriančioji plokštuma tarp skirtingų klasių. Populiariausi tiesiniai klasifikatoriai: atraminių vektorių klasifikatorius (angl. *Support Vector Machines – SVM*) ir neuroninis tinklas.[7].

*Sprendimų medžiais* grindžiamas apmokymui skirtų duomenų dekomponavimas pagal nustatytą pasirinkto atributo reikšmę [6]. Sprendimų medžio algoritmo rezultatą galima pavaizduoti struktūra, panašia į medį, kurio kiekvienas išsišakojimas reiškia vienos ar kitos sąlygos tenkinimą. Dalijimo sąlyga – žodžio arba frazės faktas. Taip sudaromos taisyklės, kurios leidžia nagrinėjamą duomenų aibę suklasifikuoti, atsižvelgiant į požymių savybes [27] Specifiniai tekstų klasifikavimui skirtos sprendimų medžių tipai – regresijos medžiai ir automatizuotas chi kvadrato sąveikos nustatymas (angl. *Chi Square Automatic Interaction Detection – CHAID*) [13].

*Taisyklėmis pagrįsti klasifikatoriai.* Taisyklėmis grindžiami klasifikatoriai modeliuoja duomenų erdvę atsižvelgdami į nustatytas taisykles. Sąlyga – nustatytų požymių, indikuojančių atitinkamą kategoriją, rinkinys.

### 1.3.2.2. Iš dalies prižiūrimas ir neprižiūrimas mokymas

Kartais sudėtinga sužymėti didelę apmokymui skirtų duomenų aibę. Todėl naudojami iš dalies prižiūrimojo ir neprižiūrimojo mokymo metodai. Šiuo atveju klasifikavimui taikomi semantinės krypties metodai (angl. *Pointwise Mutual Orientation – PMI*). Tuo tikslu naudojamas atskiras kategorijas aprašančių raktažodžių sąrašas, kuriuo remiantis nustatomas analizuojamų žodžių ir atskirų emocinių klasių raktažodžių pasiskirstymo panašumas [21]. Emocijų klasifikavimo užduotyse raktažodžiai paprastai išgaunami iš sentimentų žodyno [19].

### 1.3.3. Sentimentų žodynu grindžiamų metodų apžvalga

Išskiriami du žodynu grindžiamų metodų tipai:

- Žodyno metodas, pagrįstas sinonimų paieška žodyno generavimui nustatytiems raktažodžiams. Šis metodas aprašytas skyrelyje.
- Tekstyno metodas, kuriuo generuojami žodžiai iš apmokyto tekstyno taikant statistinius arba semantinės krypties metodus.

Tekstynas – tai sužymėtų tekstinių dokumentų rinkinys, naudojamas teksto klasifikavimui ir sentimentų žodyno generavimui.

Yra įvairių sentimentinės analizės krypčių tekstynų. WaCky [24] yra didžiulis iš žiniatinklio surinktas ir lingvistinėmis priemonėmis apdorotas tekstynas. Jame pateikiamos trijų skirtingų sluoksnių anotacijos:

1. sakinio lygmuo (objektyvus/subjektyvus, teigiamas/neigiamas/neutralus)
2. frazių lygmuo (anotuotas poliariškumas ir modifikatoriai)
3. teiginių lygmuo (suanotuotos emocinės būsenos panašiai kaip MPQA tekстыne) [9]

Pastaruoju metu publikuotas amazon.com atsiliepimų tekstynas USAGE, skirtas aspektinei analizei. Jame pateikiama po 800 vokiškų ir anglišku sužymėtų klientų atsiliepimų pagal vertinimo aspektą [20].

Be aukščiau paminėtų, yra ir daugiau iš klientų atsiliepimų sudarytų tekstynų. Kadangi atsiliepimų puslapiuose be tekstinių vertinimų, pateikiamas ir skalių vertinimas skaitinėmis išraiškomis arba žvaigždutėmis, tai tokius atsiliepimus galima interpretuoti tiesiogiai ir jie reikalauja minimalaus pradinių duomenų paruošimo [7]

## 2. Emocijų analizės eksperimentas

Emocijų analizės eksperimento tikslas – teksto emocijų klasifikavimo uždavinio įgyvendinimas. Teksto klasifikavimas buvo atliekamas taikant semantinį, mašininį-statistinį ir hibridinį metodus.

Teksto klasifikavimo uždavinio pradiniam etape buvo atliekama duomenų gavyba ir pirminis paruošimas. Surinkti duomenys sukaupti sužymėto teksto ir emocijų žodyno pavidalu. Sekančiuose skyreliuose pateikiama duomenų gavybos ir pirminio paruošimo eiga.

### 2.1. Duomenų gavyba

Emocijų ir kitų kryptų skaitmeninio teksto klasifikavimo eksperimentų sėkmės pagrindas – tinkamai atrinkti ir paruošti duomenys.

Duomenys renkami tiek iš oficialių žiniatinklio dokumentų (naujienų portalai, literatūros kūriniai), tiek iš socialinio turinio dokumentų: skaitytojų komentarų, atsiliepimų. Socialinio turinio dokumentuose aptinkama nemažai nenorminių žodžių: barbarizmų (fainas), prailgintų žodžių, kuriems būdinga balsės pakartojimai (nuooostabu), emocijų reiškimo simbolių (:D), piktogramų (2rys), gramatinių klaidų. Tokie žodžiai neatpažįstami nei morfologinio analizatoriaus, nei afektyvių terminų žodyno, todėl prieš analizuojant būtinas pirminis duomenų paruošimas. Šio etapo metu duomenys sunorminami, išvalomi nuo nereikšmingų analizei, necenzūrinių ar jautrių duomenų. Iš apdorotų duomenų konstruojami tekstynai.

Siekiant teksto emocijų klasifikavimo sistemos universalumo, duomenys buvo išgaunami iš įvairių tipų žiniatinklio dokumentų: elektroninių naujienų portalų "Lietuvos rytas" ir "Delfi" pranešimų, skaitytojų komentarų bei diskusijų forumų atsiliepimų. Surinktų duomenų pagrindu sugeneruoti du tekstynai iš skirtingų žanrų tekstų: apmokymui – iš buitinio stiliaus tekstinės informacijos (skaitytojų komentarai ir atsiliepimai), o testavimui skirti duomenys sugeneruoti iš bendrinės lietuvių kalba parašytų elektroninių naujienų portalų pranešimų ir antraščių. Nors antraštės – trumpi lakoniški tekstai, bet jais siekiama sužadinti emocijas ir patraukti potencialų skaitytoją, todėl tai labai tinkamas emocijų tekste aptikimui žanras.

Internetinė žiniasklaida – jauna žiniasklaidos rūšis, tačiau ji netruko įnešti pokyčių į žiniasklaidos tekstų organizavimo procesą, ypatingą dėmesį skirdama antraštėms. Dėl informacijos gausos, neįmanoma perskaityti visų naujienų, todėl antraštė turi patraukti potencialų skaitytoją. Naujienų portalų antraštės yra tas trumpojo teksto žanras, kuris pastebimai evoliucionuoja Suzan Kavanoz [6], gilindamasi į antraščių ir jose vartojamos kalbos paskirtį, pastebėjo, kad antraščių poveikis, palyginti su jomis įvardijamų tekstų poveikiu, skaitytojui, tikėtina, yra stipresnis. „Antraštė – svarbi žinia, kartais dargi paryškinta, išdidinta ar kitaip išskirta, pirmoji krenta į skaitytojo akis. Autoriai gerai žino, kad skaitytojas pasiduoda pirmojo akimirksnio padarytam poveikiui, todėl siekdamas atkreipti jo dėmesį nesitenkina neutraliąja antrašte, dažnai nevengiama hiperbolizuoti [6].

Teksto duomenų išgavimui įprasta taikyti žiniatinklio gavybos metodus.

Žiniatinklio gavyba – tai automatizuotas informacijos iš žiniatinklio dokumentų išgavimas. Python kalba yra parašyta keletas bibliotekų, palengvinančių žiniatinklio duomenų gavybą:

- Requests – biblioteka, skirta žiniatinklio puslapio atsisiuntimui;
- Python standartinės bibliotekos HTML analizatorius – html5 formato žiniatinklio elementų analizavimui;

- Webscraping, PyQuery, BeautifulSoup, lxml – bibliotekos, skirtos lxml/html5 formato žiniatinklio elementų analizavimui;
- Mechanize, Scrapy – bibliotekos elementų atrankai, paremtai XPath selektoriams.

Duomenys teksto emocijų klasifikavimo eksperimentui surinkti taikant žiniatinklio gavybos metodus, pasitelkiant standartinės Python bibliotekos HTML analizatorių ir reguliariąsias išraiškas. Iš žiniatinklio išgaunami duomenys buvo kaupiami SQLite duomenų bazėje. Kiekvienam įrašui priskirti šie atributai: tekstas, teksto santrauka, naudojant maišos funkciją, šaltinis, išgavimo data.

Sugeneruoti apmokomi tekstynai sužymėti pagal 6 kategorijas, atitinkančias emocijų Ekmano modelį. Sužymėta 500 apmokymui ir 200 testavimui skirtų įrašų.

Tekstyno įrašai su klasės žyma ir emocijos intensyvumo įverčiu nukopijuoti ir išsaugoti tekstinio dokumento pavidalu. Šis formatas patogesnis saityno paslauga perduodant duomenis morfologinio žymėjimo sistemai.

Tekstyno dokumentai buvo anotuojami trijų nepriklausomų anotuotojų. Testavimui buvo palikti tik tie tekstai, kurių emocinis įvertinimas bent dviejų anotuotojų buvo identiškas, o trečiojo anotuotojo įvertinimas nesiskyrė savo poliariškumu, pvz., jei pirmieji du priskyė tekstui emociją „laimė“, tai trečiojo priskirta emocija galėjo būti nuostaba, bet ne liūdesys, baimė, pyktis ar pasibjaurėjimas.

Atlikti teksto emocinio žymėjimo eksperimentai [24] atskleidė, kad priskirti tekstui emocijos klasė yra sudėtinga, ypač problematiškas nuostabos klasės nustatymas. Siekiant teisingai nustatyti emocijas tekste būtina teisinga emocijos samprata, todėl internetinių tekstų emocijų priskyrimo anotuotojai buvo supažindinti su Lazarus kognityviaja emocijų vertinimo teorija[3]. Remiantis šia teorija svarbu įvertinti faktą, sužadinusį emociją, reakcijos į tą faktą pobūdį (aktyvi, pasyvi, kuriančioji, destruktivi, teigiama, neigiama) ir šaltinį, kuris gali būti vidinis, kai emocija pergyvenama viduje ir išorinis, kai emocija nukreipta į kitą objektą, t.y. išorę. Žemiau pateikiamas pirminių emocijų įvertinimas kognityviosios teorijos požiūriu.

1. Džiaugsmas – aktyvi, kuriančioji, teigiama, nukreipta tiek į išorę, tiek į vidų emocija.
2. Nuostaba – pasyvi, vidinė emocija.
3. Liūdesys – pasyvi, destruktivi, vidinė emocija.
4. Pyktis – aktyvi, destruktivi, išorinė.
5. Baimė – pasyvi, destruktivi, vidinė.
6. Pasibjaurėjimas – pasyvi, vidinė, neigiama emocija, kuriai būdingi fiziologiniai šleikštulio požymiai.

Kita problema, su kuria susidūrė anotuotojai – teksto emocijos daugiareikšmiškumas. Tiek tekstui, tiek realioms situacijoms būdingas emocijų persipynimas, kada vienu metu išgyvenama keletas emocijų: įprasta ir džiaugtis, ir nerimauti arba pykti ir bijoti tuo pačiu metu. Tai patvirtina ir masinio atskirų žodžių anotavimo būdu (angl. crowdsourcing) sudaryto emocijų žodyno NRC rezultatai: daugumai žodžių priskirtos dvi-trys emocijos klasės. Siekiant išvengti klasifikavimo klaidos dėl teksto priklausymo kelioms klasėms, anotuotojai buvo supažindinti su emocijų raktažodžių įverčiais. Taigi tekstas buvo anotuojamas remiantis ne tik subjektyviu teksto emocijos įvertinimu, bet ir objektyviais atskirų žodžių įverčiais, nustatytais pagal dimensinę teoriją.

## **2.2. Emocijų žodyno generavimas**

Eksperimento metu emocijų žodynas buvo generuojamas naudojant hibridinį metodą. Pirminis žodynas sudarytas rankiniu būdu pagal taksonomiją, paremtą emocijų klasifikavimo teorijų taikymu.

Emocijų vertinimo taksonomija sudaryta Ekman diskrečiosios emocijų klasifikavimo teorijos, plačiai taikomos emocijų analizės uždaviniuose [3,17,4], pagrindu. Išskirtos šešios pirminės emocijos: laimė, nuostaba, pyktis, liūdesys, baimė ir pasibjaurėjimas. Taip pat atsižvelgta ir į kitų teorijų idėjas. Dimensinės teorijos principai pritaikyti matuojant emocijas keliamatėje erdveje, priskiriant emocijos klasę indikuojančiam žodžiui papildomus junglumo ir aktyvumo atributus. Kiekvienas iš atributų padidina žodžio emocijos klasės įvertį 1 balu. Bazinis kiekvieno emocijos klasės raktažodžio svoris – 1 balas. Tad, jei emocijos raktažodžiui būdingos junglumo ir aktyvumo savybės, jo svoris bus 3 balai. Remtasi ir vertinamosios teorijos subjektyvumo idėja, pirminių emocijų modelį papildant antrinėmis ir tretinėmis emocijomis.

Į pradinį emocinį žodyną įtrauktos ne tik emocijų būsenos, bet ir emocijas sužadinantys dirgikliai bei procesai, apibūdinantys emocijų raišką.

Eksperimento metu emocijų žodynas buvo generuojamas naudojant hibridinį metodą. Pirminis rankiniu būdu sudarytas afektinis žodynas buvo plečiamas žodyno ir tekstyno metodais. Žodyno metodu buvo generuojami sinonimai iš populiaraus dėl savo išsamumo ir apimties „WordNet“ žodyno, tekstyno metodu iš sužymėto tekstyno buvo išgaunami susiję su emocijų raktažodžiais žodžiai ir žodžių junginiai.

Emocijų žodyno „branduolys“ sudarytas remiantis nustatyta emocijų vertinimo taksonomija, kurios pagrindas – Ekman emocinis-psichologinis modelis [12].

Pavieniai žodžiai (unigramos) buvo plečiami iki žodžių junginių (frazių), sudarytų iš dviejų-keturių žodžių. Iš sužymėto tekstyno buvo generuojami dviejų tipų žodžių junginiai: n-gramos ir kolokacijos.

N-grama – tai gretimų žodžių seka.

Kolokacija – tai žodžių, nebūtinai gretimų, seka, aptinkama dažniau nei atsitiktiniai žodžiai [6]. Pasirinktas kolokacijų atrankos kriterijus – trys ir daugiau pasikartojimai aqnalizuojamame tekстыne. N-gramų atrankos kriterijumi pasirinktas ne n-gramos pasikartojimas, o struktūra: į emocijos požymių vektorių įtrauktos n-gramos, turinčios išreikštą emocinę prasmę ir tos, į kurių sudėtį įeina kitų emocijų klasių, nesutampančių su n-gramos klase, žodžiai. Pavyzdžiui, frazė „laimėtojas apstulbo“ priskirta nuostabos emocijos klasei, o ją sudarantis žodis „laimėtojas“ – laimės klasei. Kadangi algoritmo vykdymo metu aptikta n-grama, priskyrus jos klasės įvertį, panaikinama analizuojamame sakinyje, ją sudarantys žodžiai toliau nebeanalizuojami ir tokiu būdu išvengiama neteisingo n-gramą sudarančio žodžio klasifikavimo fakto.

N-gramų savybės plačiai taikomos daugelyje teksto analizės uždavinių, įskaitant teksto emocijų klasifikavimą [2], [16]. Šiame tyrime atlikti eksperimentai su unigramomis (n=1), bigramomis (n=2), trigramomis (n=3) ir jų deriniais patvirtino, kad žodžių junginių naudojimas sentimentinėje analizėje duoda tikslesnius rezultatus (šiuo atveju tikslumas padidėjo maždaug 5 procentais). Skyrybos ženklai taip pat įtraukti į žodžių junginių modelį. Be to dvigubi skyrybos ženklai buvo susieti su emocijne žyme, pavyzdžiui, “??” ir “?!” buvo priskirti nuostabos kategorijai.

Emocijų žodynas papildytas iš konceptualių metaforų sugeneruotais žodžiais ir frazėmis. Konceptualių, arba kitaip vadinamų abstrakčių, metaforų teorija grindžiamas nesąmoningas metaforų sudarymo mechanizmas kasdienėje kalboje ir mąstyme (Lakoff and Johnson (1980)). Kadangi emocijos yra visiškai nestruktūrizuoti konceptai, emocijų metaforos yra tipiškasis pavyzdys, pagrindžiantis pagrindinę kognityvinės lingvistikos idėją – abstrakčių konceptų suvokimą per konkrečius [struktūrizuotus] konceptus. Pavyzdžiui, emocija – abstraktus objektas – išreiškiama per materialius, apčiuopiamus objektus (skystis, indas, gyvūnas ir pan.). Konceptualios metaforos ypač reikšmingos

emocijų raiškos procesus apibrėžiančių žodžių generavimui.

Ilgą laiką buvo laikomasi nuomonės, kad emocijų metaforos turi būti universalios t.y. nepriklausančios nuo kultūrinio aspekto, kadangi emocijų kilmė ir raiška yra susijusi su fiziologiniais visų kultūrų žmonėms būdingais procesais. Tačiau Kovesces, atlikęs anglišų ir vengriškų konceptualių metaforų analizę, nustatė, kad metaforos gali būti dviejų tipų: universalios ir specifinės – būdingos atskiroms kalboms ir kultūroms [7] Šie faktai patvirtinti Yu's, Liu ir Zhao, Chen's atliktos anglų ir kinų bei Mashak anglų ir persų emocijų metaforų tyrimų rezultatais.

Emocijų metaforų tyrimų lietuvių kalbai neužfiksuota, todėl buvo sudaryta unikali konceptualių metaforų antologija, pritaikyta lietuvių kalbai ir kultūrai. Konceptualių emocijų metaforų atskaitos tašku buvo laikoma bet kurios emocijos išraiška šiais konceptais:

Emocija – indas

Emocija – Stichija

Emocija – Spalva

Emocija – Gyvūnas

Emocija – Substancija

Emocija – Vandens telkinys (laimė – jūra iki kelių, laimės vandenynas; liūdesio liūnas)

Emocija – Judėjimas (šokti/šokinėti iš laimės; panirti į depresiją )

Emocija – Objektas

Remiantis šiais konceptais iš LDT buvo generuojamos semantinės asociacijos – netiesioginiai emocijų indikatoriai, žymintys abstrakčios metaforos, iš kurios jie kildinami, srities konceptus. Pavyzdžiui, iš abstrakčios metaforos PYKTIS YRA ĮKATINTAS SKYSTIS SLEGIAMAME INDE nustatyti konceptai KARŠTIS, ĮKAISTI, SKYSTIS IR INDAS, o iš šių konceptų asociacijų išgauti žodžiai „virti, degti, sprogti, karštis, kaitinti, garuoti, lietus, veržtis“. Asociacijų generavimui buvo taikytas semantinės krypties (angl. *Pointwise Mutual Information – PMI*) metodas, kuriuo iš tekstyno buvo generuojami žodžiai, susiję su emocijos konceptais. Šis metodas, nereikalaujantis nei duomenų paruošimo, nei tekstyno įrašų emocijos klasės žymėjimo, pasirinktas dėl paprastumo ir efektyvumo [13]. Atrinktos ne mažiau kaip penkis kartus užfiksuotos asociacijos.

Dauguma su emocijų konceptais susijusių žodžių naudojami ir neemocinio pobūdžio kontekstuose, pvz., „įkaisti“ galima ir nuo emocijų, ir nuo židinio, ir nuo saulės, o „užvirti“ gali tiek iš pykčio kraujas, tiek arbata. Siekiant išspręsti daugiareikšmiškumo problemą, į emocijų žodyną įtraukti tik aiškia emocinę prasmę turintys pavieniai žodžiai, o kitų žodžių pagrindu sugeneruotos n-gramos (dvigramos ir trigramos).

Afektinių žodžių sąrašo išplėtimui pasirinktas žodyno metodas, kuris paremtas euristika, kad afektinio žodžio sinonimai indikuoja tą pačią arba artimą emociją. Kadangi WordNet šaltinyje pateikiami šie sąryšiai, jis dažnai naudojamas įvairių afektinių žodynų generavimui [16]. NLTK bibliotekoje prieinamas Vu ir Palmerio algoritmas, apskaičiuojantis sinonimų poliariškumą. Pritaikius šį algoritmą, į afektinį žodyną įtraukti tik tie sinonimai, kurių poliariškumo atstumas ne mažesnis kaip 0.5.



```

function BUILDEMOTIONLEXICON(emseeds) returns emlex
1  emlex←emseeds
2  Until done
3    if atstumas >=0.5
4      then emlex←emlex + FINDSIMILARWORDS(emlex)
5      emlex←POSTPROCESS(emlex)

```

### 1 algoritmas. Žodyno generavimo algoritmas

Pateiktame algoritmo pavyzdyje matyti, kad pirminiams emocijų klases indikuojantiems žodžiams iš WordNet išgaunami sinonimai ir į emocijų žodyną įtraukiami tik tie, kurių atstumo koeficientas, žymintis žodžių panašumą, ne mažesnis kaip 0.5.

Pagrindinė problema, su kuria susiduria sentimentinės analizės atstovai, sudarantys žodynus ne anglų kalba – mašininio vertimo tikslumas. Emocijas indikuojančių žodžių vertimo problema itin aktuali dėl emocijų nematerialumo, todėl siekiant tikslesnio lietuviškų žodžių atitikimo angliškajai versijai, buvo pasirinktas ne mašininis, bet eksperto vertimas.

## 2.2. Teksto emocijų klasifikavimas

Teksto emocijų klasifikavimo uždavinys įgyvendintas taikant šias metodikas:

- Semantinį metodą, grindžiamą emocijų žodyno ir lingvistinių taisyklių modeliu;
- Statistinį metodą, grindžiamą mašininio mokymu;

Analizuoti skirtingų tipų dokumentai: naujienų pranešimo žinutė, antraštė, skaitytojo komentaras ir kliento atsiliepimas. Skirtingų tipų dokumentai pasirinkti siekiant emocijų klasifikavimo sistemos universalumo. Dokumentai buvo analizuojami sakinio lygmenyje. Sakiniu sutartinai laikyti ir dokumentai, sudaryti iš vieno sakinio, ir iš keletos susijusių sakinių, ir antraštė.

Semantiniu metodu teksto emocija nustatyta sumuojant iš anksto nustatytus svorius, patikslinant juos atsižvelgiant į lingvistines taisykles. Mašininio regresijos metodu svoriai pritaikomi („išmokstami“) remiantis apmokomu tekstynu, sužymėtu pagal emocijų klases ir emocijos intensyvumą.

### 2.3.1. Semantinio metodo taikymas emocijų nustatymo uždavinyje

Semantinio metodo įgyvendinimui buvo taikomas taisyklėmis pagrįstas modelis, kurio pagrindas – sugeneruotas emocijų klases indikuojančių žodžių sąrašas.

Sudarytas klasifikavimo modelis paremtas semantinių ir emocijų žodynu grindžiamų klasifikavimo požymių deriniu. Emocijų žodynas sudarytas remiantis emocijų vertinimo taksonomija, konceptualių metaforų ontologija ir semantinės krypties metodu sugeneruotu n-gramų rinkiniu. Kadangi semantinis metodas nereikalauja apmokymo, apmokymui ir testavimui skirti tekstynai analizuoti bendra tvarka. Prieš atliekant teksto emocijų klasifikavimo eksperimentą, analizuojamas tekstas turi būti tinkamai paruoštas.

### 2.3.1.1. Pirminis duomenų paruošimas

Tekstas – tai nestructūrizuoti duomenys, kuriuos norint analizuoti, būtina tinkamai paruošti ir struktūrizuoti, ypač kruopštaus pirminio paruošimo reikalauja nebendrine kalba parašyti interneto lankytojų komentarai ir atsiliepimai. Kadangi naujienų portalų komentatoriai nevengia užgaulios kritikos ir necenzūrinių žodžių, komentarai buvo cenzūruojami, valomi nuo nereikšmingų analizei, nepasižyminčių emociu turiniu duomenų, panaikinamos gramatinės klaidos, privatūs duomenys (asmenvardžiai, partijų pavadinimai) koduojami.

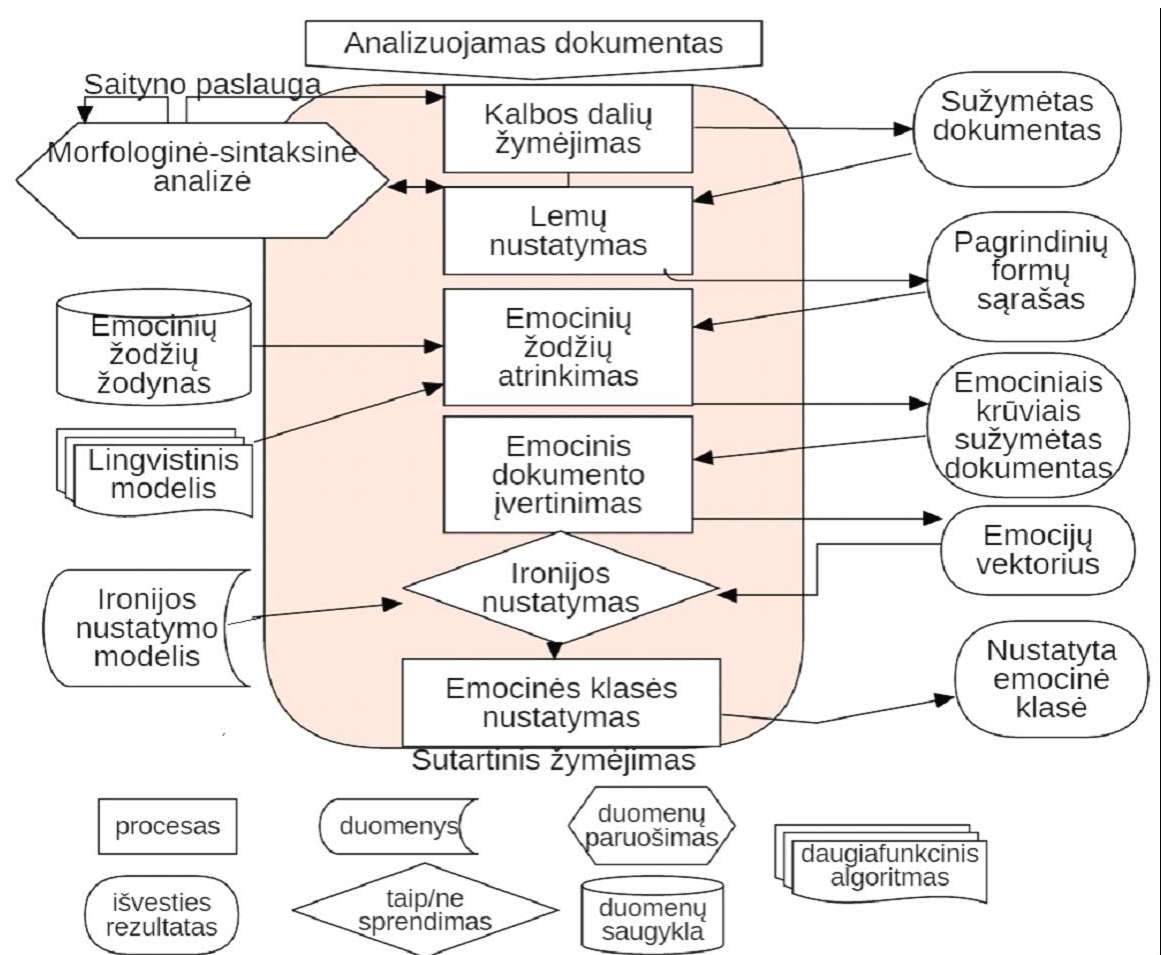
Natūralios kalbos apdorojimo priemonės buvo derinamos su žodynu grindžiama metodika siekiant išgauti semantines asociacijas ir optimizuoti teksto emocijų klasifikavimo procesą. Natūralios kalbos technologijos atliko svarbų vaidmenį pirminio duomenų paruošimo etape morfologiškai analizuojant tekstą. Morfologinis anotavimas atliekamas pasitelkiant VDU Kompiuterinės Lingvistikos Centro teksto morfologinio anotavimo saityno paslaugą. Analizuojamas tekstas buvo siunčiamas į VDU morfologinį anotatorių ir sužymėtas gražinamas emocijų analizės sistemai. Morfologinė anotavimo sistema (autorius V. Zinkevičius) pateikė tekstą sudarančių žodžių antraštines formas – lemas ir pažymėjo šių žodžių morfologines savybes, nurodančias linksnį, formas (sangražinės, neig., įvardžiutinės), asmenį. Lemų pagalba atpažįstami tekстыne esantys emocijų klasę indikuojantys žodžiai su jiems priskirtais svoriais. Morfologinės pažymos, nusakančios žodžio junginio sudėtinių dalių savybes, reikalingos šablonų, paremtų morfologiniais požymiais, atpažinimą. Pavyzdžiui, ironijai būdingi aukščiausiojo laipsnio būdvardžiai ar priešingo poliariškumo būdvardžio ir daiktavardžio seka.

### 2.3.1.2. Klasifikavimui naudingų požymių vektorius

Neapdorota tekstinė informacija – tai tik nestructūrizuoti duomenys. Siekiant klasifikuoti tekstinę informaciją būtina ją transformuoti į priimtina, struktūrizuotą formą. Dažniausiai tekstinės informacija transformuojama naudojant žodžių krepšelio (angl. *bag of words*) metodą. Žodžių krepšelio metodas – tai paprastas, tačiau efektyvus būdas atvaizduoti tekstinį dokumentą kaip jį sudarančių atrinktų požymių svorių vektorius:  $d = (f_1, f_2, \dots, f_k)$ , kur svoris  $f_k$  nusako, kiek  $k$ -asis prisideda prie dokumento  $d$  semantinės prasmės [8]. Atlikus duomenų normalizavimą ir morfologinį anotavimą, analizuojamas dokumentas pateikiamas klasifikatoriui požymių vektoriaus pavidalu. Požymių vektorius sudarytas tiek iš įprastų teksto klasifikavimo, tiek iš socialinių tinklų turiniui būdingų požymių skaitinių svorių:

- 1) Emociją indikuojančių afektinių žodžių;
- 2) Emociją indikuojančių ngramų;
- 3) Ironijos fakto;
- 4) Emocijos reiškimo simbolių;
- 5) Skyrybos ženklų;
- 6) Emocijos paneigimo;
- 7) Emocijos sustiprinimo-susilpninimo;
- 8) Emocijos poliariškumo pokyčio;
- 9) Žodžio padėties kodavimo.

Emociją indikuojantys pavieniai žodžiai ir žodžių junginiai (ngramos), skyrybos ir emocijų reiškimo simboliai paduodami sistemai masyviųjų sąrašų pavidalu, emocijos paneigimo, stiprinimo-silpninimo, poliariškumo pokyčio bei žodžio padėties kodavimas vykdomas įgyvendinant lingvistinių taisyklių modelį.



**1 pav.** Semantinio metodo schema

### 2.3.1.3. Lingvistinių taisyklių modeliavimas

Sugeneruotame emocijų žodyne nustatyti žodžio aprioriniai emocijų klasių įverčiai, priklausantys nuo analizuojamu terminu išreiškiamos emocijos intensyvumo ir junglumo. Tačiau skirtinguose kontekstuose žodžio emocinis svoris gali pasikeisti priklausomai nuo įvairių konteksto požymių: emocijos neigimo, stiprinimo, silpninimo. Konteksto požymiai apibrėžiami sintaksinėmis konstrukcijomis, indikuojančiomis analizuojamos emocijos buvimą. Todėl buvo sumodeliuotos žemiau pateikiamos lingvistinės taisyklės, patikslinančios apriorinius emocijų klasių svorius.

*Negatyvumo taisyklė:* Emocijos arba, plačiau prasme, sentimentu paneigimas sentimentinės

analizės uždaviniuose yra dažniausiai modeliuojama taisyklė [4]. Neiginių naudojimas gali pakeisti emocinio žodžio prasmę, t.y. žodžiui priskirti priešingo poliariškumo emocinę klasę arba sumažinti tos pačios emocijų klasės įverčius sakinyje ar sakinio dalyje (jei tai sudėtinis sakinytis). Į neiginių sąrašą įtraukti neigiami prieveiksmai, neigiamos dalelytės ir veiksmažodžiai, indikuojantys praradimą.

Išsamių emocijų klasių paneigimo tyrimų neužfiksuota, tačiau, kadangi emocijų nustatymas yra laikoma sentimentinės analizės kryptimi, remiamasi sentimentų paneigimo tyrimų rezultatais [8].

Nustatyta, kad labai teigiamų žodžių paneigimo atveju, keičiasi tiek poliariškumas, tiek emocijos intensyvumas – gaunamos silpnai neigiamos frazės. Tačiau labai neigiamų žodžių paneigimo atveju pakinta tik sentimentų intensyvumas, o neigiamas poliariškumas išlieka. Kiritchenko et al. atliko išsamų paneigiamų sentimentų tyrimą, kurio metu naudojant statistinį tekstyno metodą įvertinti paneigiamų individualių žodžių svorių pokyčiai, sugeneruoti atskiri šių žodžių teigiamo ir neigiamo konteksto žodynai. Remiantis eksperimento duomenimis paneigiamų emocijų raktažodžių svoriai sumažėja apie 50 procentų. Atsižvelgiant į išryškėjusias tendencijas paneigiamoms neigiamoms emocijoms išlaikoma ta pati emocinė klasė, dvigubai sumažinant jų intensyvumą. Nors teigiamoms emocijų klasėms (laimė, nuostaba) būdingas poliariškumo pokytis, tačiau vienareikšmiškai jį nustatyti sudėtinga, kadangi paneigiama laimės emocija gali pakisti į liūdesį, pyktį ar pasibjaurėjimą. Todėl teigiamų emocijų paneigimas modeliuojamas panaikinant apriorinį emocijų svorį.

Kita problema – sakinio dalies, kurioje ieškoma negatyvumo požymių, ilgis. Žvalgomojo tyrimo metu analizuojant tekstyno sakinius su neiginiais, įvertinta, kad pakanka negatyvumo žymiklių ieškoti lange, sudarytame iš šešių žodžių. Tačiau didžiojoje dalyje (86 proc.) įrašų neiginiai aptikti nutolę nuo emocinio žodžio per 1-2 pozicijas. Todėl negatyvumo taisyklė modeliuojama prieš tai suskaidant sakinį dalimis ir analizuojant kiekvieną dalį atskirai. Sakinio dalijimo kriterijais pasirinkti skyrybos ženklai ir prieštaros jungtukai.

*Prieštaros taisyklė:* Neretai vienareikšmiškai identifikuoti emociją gali būti sudėtinga dėl viename sakinyje (dokumente) aptinkamų skirtingas emocijas indikuojančių žodžių. Tai ypač būdinga sudėtiniais priešinamiesiems sakiniams

Modeliuojant emocinio klasifikavimo algoritmą ypač reikšmingi sudėtiniai sakiniai su prieštaros semantikos jungtukais „bet“, „tačiau“, „o“. (pvz., „Nors jis ir turi trūkumų, bet iš esmės yra geras žmogus“). Priešinamuosiuose sakiniuose nuomonė prieš prieštaros jungtuką ir po jo paprastai prieštarauja viena kitai. Tai taikoma ir kai kurioms kitoms frazėms, pvz. su žodžiais „išskyrus“, „nepaisant“, todėl prieštaros jungtukų sąrašas sudarytas atsižvelgiant į prieštaros semantiką. Modeliuojant šią taisyklę buvo koduojama žodžio padėties sakinyje informacija. Remtasi hipoteze, kad žodžiai, išdėstyti sakinio gale, turi didesnę emocinį svorį, kadangi žmonėms būdinga apibendrinti ar išryškinti savo požiūrį sakinio gale. Todėl padėties informacija buvo koduojama, kiekvienam terminui priskiriant atributus, žyminčius žodžio padėtį dokumente (iki ir po prieštaros jungtuko). Šiems atributams skaičiuojant emocijų žodžių įverčius priskiriami koeficientai: 0.5 – iki ir 2 – po. Tokiu būdu emocinio žodžio svoris apskaičiuotas dauginant emocijų žodžių žodyne nustatytą žodžio įvertį iš koeficiento:

$$w_e = w_e * PI, \text{ kur } PI - \text{ koeficientas, } w_e - \text{ emocinio žodžio } w \text{ svoris} \quad (3)$$

*Stiprinančiųjų prieveiksnių taisyklė:* Dvigubas svoris buvo skiriamas ir emocijoms žodžiams, naudojamiems kartu su stiprinančiosiomis išraiškoms, pvz. geras išraiškoje „labai geras“ įgis dvigubai didesnę vertę nei „geras“ kontekste be „labai“ ar kitų stiprinančiųjų žodžių ar frazių.

*Padidrinimo-sumažinimo taisyklė:* Ši taisyklė teigia, kad savybės, išreikštos emociiniu žodžiu, stiprumas ( kiekis, poliškumas) gali keistis priklausomai nuo greta analizuojamojo žodžio esančių žodžių semantikos. Veiksmazodžiai gali turėti emocijos padidrinimo (igijimo) arba sumažinimo (trūkumo) prasmę. Pavyzdžiui, „Vaistai sumažina skausmą”. Nors „skausmas” yra neigiamo poliariškumo liūdesio emocijos klasės žodis, tačiau sumažinimas šiuo atveju reiškia pageidaujamą teigiamą laimės emociją atitinkantį rezultatą. Pavadinimas „sumažinimo-padidrinimo“ sąlyginis, kadangi į sąrašus įtraukti ir emocijos atsiradimo bei netekimo veiksmazodžiai. Nors emocijos praradimo semantika artima emocijos paneigimo semantikai, tačiau į negatyvumo sąrašą įtraukti neiginiai, sudaryti iš pagalbinių kalbos dalių (dalelytė, prieveiksmai), o sumažinimo (panaikinimo) sąrašas sudarytas iš veiksmazodžių. Taip pat skiriasi žymiklių paieškos sritis: negatyvumo požymių ieškoma visoje semantiškai savarankiškoje sakinio dalyje, o panaikinimo – analizuojami tik žodžiai, nutolę per 1-2 pozicijas. Tokiems atvejams naudojami didrinimo ir mažinimo prasių veiksmazodžių sąrašai ir skaičiuojama pagal šias taisykles :

veiksmazodis + liūdesio/baimės/pykčio/pasibjaurėjimo emocinis žodis → tos pačios klasės emocija su 2/3 sumažintu liūdesio emocijos svoriu, pvz. sakinyje „Mano problema visiškai išnyko” *problema*, emocinių žodžių žodyne pažymėta baimės emocija su svoriu 1, įgis laimės emocijos klasę su įverčiu 0.5. Be to, remiantis stiprinančiųjų prieveiksnių taisykle stiprinančiojo žodžio „visiškai“ dėka emocinis svoris padvigubės;

Didrinimo veiksmazodis + liūdesio/baimės/pykčio/pasibjaurėjimo emocinis žodis → liūdesio/baimės/pykčio/pasibjaurėjimo emocija su dvigubu svoriu;

Didrinimo veiksmazodis + laimės/nuostabos emocinis žodis → laimės/nuostabos emocija su 2/3 padidintu svoriu;

Mažinimo veiksmazodis + laimės emocinis žodis → liūdesio emocija su laimės emocijos svoriu, pvz., „Dėl darbo jis atsisakė asmeninės laimės.” laimė bus vertinama kaip žodžio „laimė“ vertės liūdesio emocija; (norėti)

Mažinimo veiksmazodis + nuostabos emocinis žodis → nuostabos emocija su dvigubai mažesniu emocinio žodžio svoriu.

6. *Trūkumo veiksmazodžių taisyklė:* Esant trūkumo veiksmazodžiams, pvz. norėti, siekti, reikėti, pakinta tame pačiame sakinyje esančių emocinių žodžių įvertis ir poliškumas. Remiantis euristika, kad žmonėms būdinga siekti teigiamų dalykų, ši taisyklė pritaikyta laimės emocinės klasės žodžiams, kuriems priskiriama liūdesio emocijos klasė su 3/4 apriorinio laimės emocijos svorio. Kaip ir negatyvumo, tik skirtingas koeficientas ir tik laimės klasei.

#### **2.3.1.4. Klasifikavimo rezultatai**

Dokumento klasifikavimas buvo grindžiamas žodžių analizuojamame dokumente palyginimu su emocijų žodyne esančiais žodžiais, turinčiais kiekvienos emocijos nustatytus įverčius (svorius). Žodžių emociniai svoriai sakinyje sumuojami pagal atskiras emocijas, dokumentui priskiriama emocija, kurios svoris didžiausias. Svoriai nustatyti pagal sudarytą emocijų vertinimo taksonomiją ir patikslinti klasifikavimo metu atsižvelgiant į sumodeliuotas lingvistines taisykles. Nustačius emocijos klasę buvo tikrinamas dokumentų su nustatyta laimės emocija ironijos faktas. Jeigu ironijos neaptikta, emocijos klasė patvirtinta, jeigu ironijos faktas nustatytas, dokumentui priskirta pykčio arba pasibjaurėjimo emocija priklausomai nuo didesnio šias emocijas indikuojančio įverčio analizuojamame dokumente.

Ironijos nustatymui sudaryta atskira posistemė, kurios aprašymas pateikiamas sekančiame skyrelyje.

### 2.3.1.5. Ironijos nustatymo algoritmas

Atliktų susijusių darbų išvadomis nustatyta, kad vienas pagrindinių sentimentinės analizės ribojančių veiksnių yra perkeltinės kalbos naudojimas [18] ištyręs ironijos suvokimo ir tikrosios reikšmės dekodavimo mechanizmus, atskleidė, kad ironija kaip retorine priemone siekiama sustiprinti reiškiamą emociją. Kitų eksperimentų išvadomis teigiama, kad ironijos faktas susilpnina išreiškiamą emociją. Tačiau tyrimais pagrįsta, kad ironijos fakto nustatymas pagerina teksto emocijų klasifikavimo kokybę [18].

Ironija kaip klaidingo klasifikavimo priežastis būdinga tik klaidingai priskirtoms teigiamų emocijų klasėms [Liu]. Kadangi pagal sudarytą emocijų vertinimo taksonomiją visiškai teigiama emocija laikoma laimė, tai ironijos nustatymas taikomas tik dokumentams, kuriems priskirta laimės emocijos klasė.

Skirtingi autoriai ironiją vertina skirtingai. Pavyzdžiui, Gibbs [4] ironiją apibrėžia kaip įvairių meninės raiškos priemonių – sarkazmo, hiperbolės, jumoro, retorinio klausimo – visumą. Kiti autoriai [4] ironiją ir sarkazmą priskiria skirtingoms meninių raiškos priemonių kategorijoms (sarkazmas - piktas pašiepimas, aštri ironija; kandi, pajuokiamą pastaba, ironija - paslėpta sąmojinga pašaipą, pasityčiojimas). Šiame darbe sudarant ironijos tekste aptikimo modelį remiamasi Gibbs ironijos samprata.

Išskiriamos 3 pagrindinės ironijos rūšys:

- situacinė ironija, žyminti neatitikimą tarp to, ko tikimasi ir kas iš tikrųjų įvyksta;
- dramatinė ironija. Ši ironijos forma panaši į situacinę ironiją, tačiau sutinkama kino, literatūros, teatro kūrinuose. Situacijos ironiškumas pasireiškia tuo, kad žiūrovas žino tai, ko nežino veikėjas;
- verbalinė ironija – sąmoningas žodžių naudojimas skirtinga (paprastai priešinga) reikšme negu jų tikroji semantinė reikšmė.

Teksto analizės užduotims aktuali tik verbalinė ironija, todėl toliau darbe sutinkamas terminas „ironija“ naudojamas verbalinės ironijos prasme. Išskiriamos šios pagrindinės verbalinės ironijos formos:

- Hiperbolė. Tai dirbtinis savybių išdidinimas (perdėjimas) ;
- Sumenkinimas (angl. understatement) – kai savybės dirbtinai sumenkinamos;
- Sarkazmas, kurį vieni autoriai laiko atskira retorine figūra, kiti – ironijos rūšimi. Nors sarkazmas pasižymi siekimu įžeisti, nebūdingu ironijai, tačiau šiame darbe sarkazmas vertinamas kaip aštresnė ironijos forma.

Ironijos formos (perdėjimas ir sumenkinimas) implikuoja ironijos raiškos priemones: perdėjimui naudojami itin stiprų teigiamą emocinį svorį turintys žodžiai ir žodžių junginiai, sumenkinimui – mažybinės formos. Kadangi žvalgomojo ironijos tekstyno vertinimo metu sumenkinimo faktų neaptikta, tai į ironijos indikatorių sąrašą įtrauktas tik hiperbolės mechanizmas. Šis ir kiti ironijos klasės nustatymo požymiai pasirinkti remiantis euristiniu metodu ir atliktų susijusių darbų išvadomis [5], jų koreliacija įvertinta chi-kvadrato testu. Sudarytas ironijos požymių vektorius pateikiamas lentelėje.

**3 lentelė.** Ironijos požymių lentelė

Ironijos klasės požymis	Požymio aprašymas
Hiperbolė	Stiprinamieji prievieksmiai+ įprasti teigiami būdvardžiai/prievieksmiai/daiktavardžiai
Žodžiai	Stipriai teigiami žodžiai
Frazės	Kalbos šablonai, būdingi ironijai saugomi ironijos žodyne
Jaustukai	oi, ai, oh, ah, omg, vau, juk, aha, vaje, vajė, oj, ot, ale, chi, hi, cha
Sintaksinės priemonės	Skyrybos ženklai (klaustukas, daugtaškis, kabutės)
Sutartiniai emocijų simboliai	;), :D, P
Pakartojimai	Beasmenės formos pritariančiųjų išraiškų pasikartojimai
Oksimoronas	Priešingo poliariškumo žodžių seka

Ironijos žodynas sugeneruotas iš rankiniu būdu sudaryto ir trijų nepriklausomų anotuotojų sužymėto tekstyno.

Kiekvienas analizuojamas dokumentas pateikiamas ironijos klasės požymių vektorius forma. Ironijos klasės požymiams priskirta loginė reikšmė: 1 – jeigu požymis būdingas analizuojamam dokumentui, 0 – jeigu požymio analizuojamame dokumente neaptikta. Ironijos faktas patvirtinama bent dviem vektorius elementams priskyrus teigiamą loginę reikšmę.

Rezultatai: Analizuota 204 dokumentai, iš jų 133 su ironijos faktu. Teisingai teigiamų – 88, klaidingai teigiamų – 21. Tikslumas –0.66, atpažintų objektų kiekis – 0.74. F-įvertis – 0.9.

### **2.3.2. Mašininio-statistinio metodo taikymas emocijų nustatymo uždavinyje**

Tačiau vien taisyklėmis ir emocijų žodžių žodynu grindžiamo teksto emocijų klasifikavimo tikslumas – apie 30 proc. (pagal F-įvertį) Emociniai žodžiai sudaro 23 proc. visų tekstyno žodžių, tačiau analizuojamai emocijų klasei gali būti reikšmingi ir į emocijų žodžių žodyną neįtraukti žodžiai. Be to, išsamaus afektinių žodžių žodyno sudarymas reikalauja daug laiko resursų, todėl vis dažniau teksto analizės uždaviniuose taikomi mašininiai ir/arba statistiniai metodai naudingų klasifikavimui požymių atrinkimui.

#### **2.3.2.1. Naudingų klasifikavimui požymių atrinkimas taikant statistinį metodą**

Naudingi klasifikavimui požymiai gali būti ne tik žodžiai, bet ir su emocijų klase susiję simboliai. Todėl į emocijų klasifikavimui naudingų požymių atrinkimo modelį įtraukti ne tik tekstyną sudarantys žodžiai, bet ir simboliai, kognityvinės lingvistikos eksperimentų išvadamis [11] nustatyti reikšmingais tam tikrų emocijų identifikavimui.

Tokie žodžiai gali būti išgaunami semantinės krypties metodu, kuriam reikalingas pirminis emocijų indikuojančių raktažodžių sąrašas. pvz. Turney išrastas asociacijų metodas.

Tačiau paprasčiau reikšmingų požymių nustatymui taikyti statistinius metodus. Remiantis atliktų

tyrimų duomenimis [11] aukštesnis klasifikavimo tikslumas stebimas teksto klasifikavimo užduotyse naudingų požymių atrinkimui taikant chi kvadratų testą todėl šiame eksperimente reikšmingi emocijų klasifikavimui požymiai iš apmokomo tekstyno nustatyti pasitelkiant chi-kvadrato  $\chi^2$  testą, kuriuo kiekvienam tekstyno elementui apskaičiuotas chi kvadrato kriterijus ir atrinkti elementai, kurių kriterijus prilygo 3.84 ir daugiau.

Šis testas paremtas idėja, teigiančia, kad jei žodžio aptikimo tikimybė kažkurios klasės tekste yra didesnė negu kitų emocijų klasių tekstuose, tai toks žodis arba simbolis gali būti laikomas emocijos indikatoriumi.

Statistikoje  $\chi^2$  testu siekiama nustatyti dviejų įvykių tarpusavio nepriklausomybę. Įvykiai A ir B laikomi nepriklausomais, jei  $P(AB)=P(A)P(B)$  arba analogiskai  $P(A|B)=P(A)$  ir  $P(B|A)=P(B)$ . Naudingų klasifikavimui požymių atrinkimo atveju įvykiais laikomi termino (indikatoriaus) ir klasės nustatymo faktai. Skaičiavimui naudojami emocijos klasėmis sužymėto tekstyno įrašai. Iš pradžių laikomasi nulinės hipotezės, kad tiriamasis indikatorius w nepriklauso nei vienai emociinei klasei. Pearsono chi-kvadratų testas palygina gautus w dažnius su tikėtinais dažniais. Lentelėje atsitiktinumų  $p_{ij}$  žymi indikatoriaus pasikartojimo dažnumą atskirų emocijų klasių įrašų aibėje, pvz. p11 nurodo laimės klasės įrašų skaičių, kuriuose aptiktas analizuojamas indikatorius.

**4 lentelė.** Požymių tikimybių skaičiavimas

Eil. nr.	Emocijos klasė	Dokumentų su w skaičius	Dokumentų be w skaičius	Viso emociinės klasės dokumentų
1.	laimė	p11	p12	p11+p12
2.	nuostaba	p21	p22	p21+p22
3.	liūdesys	p31	p32	p31+p32
4.	pyktis	p41	p42	p41+p42
5.	baimė	p51	p52	p51+p52
6.	pasibjaurėjimas	p61	p62	p61+p62
		p11+p21+p31+p41+p51+p61	p12+p22+p32+p42+p52+p62	

Chi-kvadratų reikšmė apskaičiuojama pagal formulę:

$$\chi_i^2 = \frac{n * F(w)^2 * (p_i(w) - P_i)^2}{F(w) * (1 - F(w)) * P_i * (1 - P_i)} \quad (4)$$

$p_i(w)$  – sąlyginė klasės i tikimybė dokumentams su žodžiu w. Tikimybė – tai statistinis įvertis, nurodantis palankių atsitiktinumų santykį su bendru atsitiktinumų skaičiumi.

$P_i$  – bendras klasės i dokumentų skaičius

$F(w)$  – bendras dokumentų su žodžiu w skaičius

n – bendras dokumentų skaičius



Kuo didesnis chi-kvadratas, tuo labiau w reikšmingas emocinei klasei. Pasirinkta emocijos klasei būdingu požymiu laikyti žodį, kurio chi-kvadrato reikšmė nemažesnė kaip 3.84. Ši reikšmė atitinka reikšmingumo koeficientą 0.05, t. y. nulinė hipotezė, kad žodis-indikatorius ir emocijos klasė yra nesusiję yra atmetama, paliekant tik 5 proc. atsitiktinio įvykių sutapimo tikimybę. Atskirai suskaičiuoti požymiai, kurių  $\chi^2 \geq 6.63$ , atitinkantį reikšmingumo koeficientą 0.01. Nulinė hipotezė, teigianti, kad emociją indikuojančio žodžio ir emocijos klasė yra nesusiję, yra atmetama, paliekant tik 1 procento klaidos tikimybę.

Pyktis	Laimė	Liūdesys	Baimė	Nuostaba
ministras::18.24	gražus::35.48	gėda::43.02	drebėti::104.16	apstulbti::53.17
visas::16.98	geras::34.53	skaudėti::29.01	šurpas::71.88	žadas::50.75
tauta::13.66	aktorė::18.58	nusivilti::28.69	krėsti::41.26	tikėti::37.18
projektas::13.66	visada::17.51	seimas::16.85	virpėti::32.76	patikėti::31.52
universitetas::13.65	protingas::17.4	gaila::14.86	epušė::30.72	netekti::27.33
interesas::13.65	abu::15.01	prarasti::10.80	baimė::30.72	keistas::26.26
velnias::13.07	elegancija::14.84	vyresnis::8.37	įbauginti::30.72	nesitikėti::19.63
reikalas::13.07	nuostabus::14.84	skausmas::8.37	lapas::30.72	nepatikėti::16.33
direktorė::9.7	gražuolė::14.84	sulaukti::8.37	koja::28.7	gal::17.57
teisėtvarka::9.09	džiaugtis::14.84	verkti::8.37	krūpčioti::28.7	tekti::15.05
baikit::9.09	šaunuolis::14.84	ponas::8.37	širdis::27.32	nesuprasti::14.96
vardan::9.09	šaunuolė::14.84	neblogas::8.37	siaubas::24.88	suprasti::13.25
pažiūrėti::8.82	super::14.84	pietūs::8.37	kūnas::24.55	nustebti::13.04
paimti::8.82		vienintelis::8.37	kelias::16.2	išvysti::12.02

**5 lentelė.** Reikšmingiausių požymių taikant chi-kvadrato algoritmą sąrašas

Pasibjaurėjimas
šlykštus::121.37
fui::43.05
dvokti::37.61
šlykštėti::32.18
šlykštynė::21.38
veikėjas::16.01
apgailėti::16.01
skandalas::16.01
šmeižtas::10.65

### 2.3.2.2. Teksto emocijų klasifikavimas taikant regresijos metodą

Atrinktų klasifikavimui reikšmingų duomenų aibė 1280 unikalių žodžių. Remiantis atliktų tyrimų išvadomis, teigiančiomis, kad mašininio mokymo algoritmai pasiekia geresnius rezultatus su mažesne duomenų aibe (<30 000 unikalių žodžių/simbolių/savybių) nutarta emocijų klasifikavimo uždavinį įgyvendinti taikant prižiūravimo mokymo regresijos metodą [31]

Regresijos tikslas – nustatyti reikšmingų klasifikavimui žodžių koeficientus. Tikslą įgyvendinimui apmokyta funkcija (žr. formulę), priskirianti analizuojamą sakinį sudarantiems žodžiams svorius taip kad jų suma būtų kuo artimesnė rankiniu būdu nustatytam dokumento emocijos skaitiniam įverčiui

$$f(\vec{x}) = w_0 + \sum_i w_i * x_i \quad (5)$$

**6 lentelė.** Regresijos funkcijos sutartinis žymėjimas

Simbolis	Reikšmė
$w_i$	analizuojamo žodžio svoris
$x_i$	žodžio koeficientas analizuojamame sakinyje
$w_0$	vidutinis analizuojamos emocijos įvertis tekстыne

Funkcija skaičiuojama iteratyviai visiems tekstyną sudarantiems dokumentams. Pirmoje iteracijoje priskiriami atsitiktiniai reikšmingų klasifikavimui žodžių svoriai. Apskaičiuota funkcija, t.y. nustatytas dokumento emocijos įvertis lyginamas su realiu dokumento įverčiu. Skaičiuojama klaida, kuri minimizuojama taikant gradientinio nusileidimo algoritimą, kuriuo kiekvienoje iteracijoje svoriai tikslinami proporcingai svorio daliai išvestinei (žr. formulę).

$$w_i \leftarrow w_i - \eta \times (\partial \text{Error}_E(w)) / (\partial w_i) \quad (6)$$

kur  $\eta$ , gradientinio nusileidimo žingsnio dydis, vadinamas mokymosi greičiu. Mokymosi greitis, žodžių svoriai ir nustatytas dokumento įvertis duoti kaip įvesties parametrai. Dalinė išvestinė nurodo pokyčio dydį, minimizuojantį klaidą. Kvadratinių klaidų suma yra visų dokumentų klaidų suma toje iteracijoje. Sumos dalinė išvestinė yra visų dalinių išvestinių suma. Tokiu būdu galima analizuoti kiekvieną dokumentą atskirai ir nustatyti svorių pokytį. Dokumento  $e$  klaida yra dalinė išvestinė svorio  $w_i$  atžvilgiu:

$$2 \times [\text{val}(e, Y) - p\text{val}^w(e, Y)] \times \text{val}(e, X_i) \quad (7)$$

Tarkime, kad kiekvienam dokumentui  $e$  klaida  $\delta = \text{val}(e, Y) - p\text{val}^w(e, Y)$ . Tada kiekvienas dokumentas  $e$  patikslina kiekvieno žodžio svorį  $w_i$  dydžiu:  $\delta = \text{val}(e, Y) - p\text{val}^w(e, Y)$ . Atnaujintas svoris:

$$w_i \leftarrow w_i + \eta \times \delta \times \text{val}(e, X_i), \quad (8)$$

konstanta 2 neįtraukta į atnaujinto svorio skaičiavimo formulę, laikant, kad ji įskaičiuota į konstantą  $\eta$ . Algoritmo realizacija ir parametrai pateikiami pseudokodo pavyzdyje .

**Procedure** LinearLearner( $X, Y, E, \eta$ )

```
1: Ivestis
2:  $X$ : įvedamų požymių (regresorių) vektorius,  $X = \{X_1, \dots, X_n\}$ 
3:  $Y$ : priklausomas parametras
4:  $E$ : apmokomo tekstyno dokumentų rinkinys
5:  $\eta$ : mokymosi koeficientas (greitis)
6:  $\delta$ : klaida
7: Išvestis
8: žodžių svoriai  $w_0, \dots, w_n$ 
9:  $w_0, \dots, w_n$ : realieji sk.
10:  $pval^w(e, Y) = w_0 + w_1 \times val(e, X_1) + \dots + w_n \times val(e, X_n)$ 
11: inicializuojame  $w_0, \dots, w_n$  atsitiktiniais skaičiais
12: repeat
13: for each dokumentui  $e$  in  $E$  do
14:      $\delta \leftarrow val(e, Y) - pval^w(e, Y)$ 
15:     for each  $i \in [0, n]$  do
16:          $w_i \leftarrow w_i + \eta \times \delta \times val(e, X_i)$ 
17: until pabaiga
18: return  $w_0, \dots, w_n$ 
```

**2 algoritmas.** Regresijos funkcijos apmokymas

Algoritmą galima laikyti inkrementiniu gradientiniu nusileidimu dėl tų pačių žodžių svorių priskyrimo kiekvienam dokumentui atskirai. Tačiau jeigu atskiroje iteracijoje gauta klaida neviršija nustatyto slenksčio (0.001), žodžių svoriai, su kuriais pasiektas klaidos minimizavimas, nebetikslinami, o išsaugomi ir kitiems dokumentams bei kitoms iteracijoms naudojami statiniai šių žodžių įverčiai.

Chi-kvadrato testu iš apmokamo tekstyno atrinkti reikšmingi klasifikavimui požymiai, kuriems regresijos metodu sugeneruoti svoriai. Iš atrinktų žodžių ir jiems priskirtų koeficientų sudarytas požymių vektorius, kuris buvo naudojamas testavimui skirtame tekстыne.

**2.3.3. Eksperimento rezultatų vertinimo kriterijai**

Atliktų eksperimentų metu buvo vertinama teksto emocijų klasifikavimo kokybė pagal pasirinktų metrikų rinkinį. Automatinio būdu priskirta teksto emocijos klasė ir jos skaitinis įvertis buvo lyginama su anotuotojų priskirtomis klasėmis ir jų įverčiais. Kiekvienai klasei skaičiuojamas bendras teisingai priskirtų objektų kiekis ir objektų su priskirta klasės žyma kiekis. Šie rodikliai, atitinkantys tikslumo (angl. *precision*) ir atpažintų objektų skaičių (angl. *recall*), įtraukti į vertinimo metrikų rinkinį. Remiantis tyrimų išvadomis [], tikslumo metrikos informatyvumas esant daugiaklasiam klasifikavimui tampa mažiau reikšmingas ir realiai įgyvendinamas. Šiuo atveju svarbu ne tik, ar teisingai identifikuota emocijos klasė, bet ir klasės emocijos intensyvumo neatitikimo rodiklis. Todėl į klasifikavimo kokybės

vertinimo kriterijų sąrašą įtraukta kvadratinės paklaidos metrika. Paklaida – tai vidutinis kvadratų skirtumas tarp rankiniu būdu priskirto ir algoritmo apskaičiuoto dokumento klasės įverčio – apskaičiuota pagal formulę:

$$\frac{1}{E} \sum_{e \in E} (\text{score}(e) - \text{prediction}(e))^2 \quad (9)$$

**7 lentelė.** Vidutinės paklaidos skaičiavimo sutartinis žymėjimas

Simbolis	Reikšmė
score(e)	Rankiniu būdu nustatytas dokumento įvertis
prediction(e)	Algoritmo apskaičiuotas dokumento įvertis
E	Bendras dokumentų skaičius

Pasirinktų metrikų rinkinys susideda iš keturių metrikų:

- tikslumas (angl. *precision*);
- klaidų kiekis (angl. *error rate*);
- atpažintų arba klasifikuotų objektų kiekis (*recall*);
- F-įvertis.

Tikslumas apibrėžiamas kaip teisingai klasifikuotų objektų ir eksperto nustatytų analizuojamos klasės objektų santykis.

Atpažintų objektų kiekis apibrėžiamas kaip teisingai klasifikuotų objektų ir visų atpažintų objektų kiekio santykis.

F įvertis paremtas tikslumo ir klasifikuotų objektų skaičiaus metrikomis ir apskaičiuojamas pagal formulę (10). Kuo rezultatas artimesnis vienetui, tuo aukštesnė klasifikavimo uždavinio kokybė.

$$F - \text{score} = 2 * \frac{\text{precision} * \text{recall}}{\text{precision} + \text{recall}} \quad (10)$$

Gauti klasifikatoriaus rezultatai pateikiami lentelėje.

**8 lentelė.** Semantinio algoritmo rezultatai

	Teisingi teigiami	Klaidingi teigiami	Realus klasės objektų sk.	Tikslumas	Atpažintų objektų sk.	F-įvertis	Klaidų kiekis
Pyktis	6	9	92	0.07	0.4	0.12	0.88
Nuostaba	18	2	116	0.16	0.9	0.26	0.74
Liūdesys	5	25	55	0.1	0.17	0.15	0.85
Pasibjaurėjimas	28	13	81	0.35	0.68	0.47	0.53
Baimė	15	12	43	0.35	0.56	0.43	0.57
Laimė	102	257	105	0.97	0.28	0.44	0.56

8 lentelėje pateikiami klasifikavimo naudojant tik pavienius žodžius. Papildžius algoritmą ngramomis rezultatai nežymiai pagerėjo. Nors rezultatai nėra aukšti, tačiau pasibjaurėjimo emocijos F-įvertis artimas 0.5, nedaug atsilieka baimės ir laimės emocijų nustatymo kokybė. Eksperimento rezultatai išryškina itin aukštą klaidingai teigiamų laimės emocijų rezultatą. Eksperimento planavimo metu tai buvo numatyta, tam sukurtas ir įgyvendintas ironijos nustatymo algoritmas. Pritaikius šį algoritmą klaidingai teigiamų laimės emocijos faktų sumažėjo 28 procentais.

## Apibendrinimas

Magistro baigiamajame darbe nagrinėjamas lietuviškų žiniatinklio dokumentų emocijų klasifikavimo uždavinys. Šis uždavinys sprendžiamas taikant natūralios kalbos ir mašininio apmokymo technologijas. Pagrindinis darbo tikslas – sukurti emocijų analizės sistemos prototipą lietuvių kalbai. Kuriamo prototipo paskirtis – suklasifikuoti tekstą remiantis Ekmano emociniu modeliu ir pateikti rezultatus:

- emocijos komponentus, pagal kuriuos tekstas priskiriamas emocinei klasei;
- atsakymą, ar tekste stebimi ironijos elementai;
- atsakymą, kuriai emocijos klasei priklauso tekstas.

Emocinė analizė atlikta remiantis semantinės krypties ir statistiniu metodais, pasitelkiant žodžių bei frazių taksonomiją, lingvistinėmis taisyklėmis grindžiamą modelį, statistiniu metodu nustatytų reikšmingų klasifikavimui požymių sąrašą, kuriems mašininio mokymo metodu apskaičiuoti koeficientai.

Taikant emocijų vertinimo taksonomiją sudaryti emocijos klasę indikuojančių žodžių (1291 žodis) ir ngramų (143 ngramos) sąrašai.

Kadangi eksperimente naudoti duomenys buvo kruopščiai atrenkami ir anotuojami ekspertų, tai išgauti šablonai, stipriai koreliuojantys su analizuojama emocija, gali būti panaudoti automatiniam, nereikalaujančiam rankinio anotavimo duomenų rinkinio generavimui.

## **Pagrindiniai rezultatai**

Vykdamas šį darbą buvo pasiekti šie rezultatai:

1. Sukurtas ir įgyvendintas semantinės krypties teksto emocijų klasifikavimo algoritmas, grindžiamas sugeneruoto emocijų žodyno ir lingvistinių taisyklių modelio taikymu. Teksto klasifikavimo emocijų klasių F-įverčio rezultatų vidurkis – 0.31.
2. Sukurtas ir įgyvendintas statistinis teksto emocijų klasifikavimo algoritmas, grindžiamas regresijos ir gradientinio nusileidimo algoritmų taikymu.
3. Taikant emocijų vertinimo taksonomiją sudaryti emocijos klasę indikuojančių žodžių (1291 žodis) ir ngramų (143 ngramos) sąrašai.
4. Semantinio metodo rezultatų optimizavimui sukurtas ir įgyvendintas ironijos algoritmas, kuriuo laimės emocijos klaidingai teigiamų objektų kiekis sumažėja apie 30 procentų. Vien ironijos nustatymo algoritmo eksperimento rezultatas – 0.9 (F-įvertis).
5. Remiantis emocijų vertinimo teorijomis ir konceptualiomis metaforomis sudarytas emocijų klases indikuojančių žodžių ir ngramų sąrašai.
6. Sumodeliuotos lingvistinės taisyklės, kurios paremtos emocijos poliariškumą ir intensyvumą modifikuojančiais žodžių sąrašais.
7. Sistemos apmokymui ir testavimui sugeneruoti skirtingų diskursų duomenų rinkiniai – naujienų ir komentarų tekstynai.

## **Išvados ir rekomendacijos**

Darbo metu buvo sukurtas įgyvendintas emocijų nustatymo internetiniame tekste algoritmas, pagrįstas emocijos klasę indikuojančių žodžių, frazių ir simbolių žodynais bei lingvistinių taisyklių modeliu.

1. Darbo metu buvo nustatyta, kad aukštesnė teksto klasifikavimo kokybė gaunama semantiniu metodu: semantiniu ir statistiniu metodu pasiekti F-įverčio klasifikavimo rezultatai atitinkamai – 0.31 ir 0.17.
2. Tyrimo metu taip pat išryškėjo netolygus klasifikavimo klaidų pasiskirstymas tarp skirtingų emocijų klasių: daugiau klaidingai teigiamų priskyrimų stebima laimės emocijos klasėje negu pykčio, baimės, liūdesio ar pasibjaurėjimo. Eksperimentu patvirtinta, kad problemos sprendimui tinkamas ironijos nustatymo algoritmo integravimas į emocijų klasifikavimo sistemą.
3. Tiek statistiniu, tiek semantiniu metodais aukščiausia klasifikavimo kokybe pasižymi pasibjaurėjimo klasės emocijos nustatymas: atitinkamai F-įvertis – 0.47 ir 0.26.

Didžiausias eksperto ir mašinos emocijų indikuojančių žodžių sutapimas stebimas taip pat pasibjaurėjimo klasėje. Didžiausias chi koeficientas – 121.37 nustatytas pasibjaurėjimo klasės žodžiui „šlykštus“.



## **Ateities tyrimų gairės**

Siekiant geresnės teksto emocijų klasifikavimo kokybės būtų galima padidinti apmokomo tekstyno duomenų aibę ir papildyti sudarytą emocijų žodyną sugeneruotais žodžiais ir frazėmis iš apmokomo tekstyno. Būtų galima apjungti semantinį ir statistinį metodus ir įgyvendinti hibridinį teksto emocijų klasifikavimo algoritmą.

## Literatūros šaltiniai

- [1] Balahur A Methods and resources for sentiment analysis in multilingual documents of different text types. Dissertation, Department of Software and Computing Systems, University of Alacant, Alacant, 2011.
- [2] Baroni M, Bernardini S, Ferraresi A, Zanchetta E The WaCky wide web: a collection of very large linguistically processed web-crawled corpora. *Lang Resour Eval* 43(3):209–226, 2009.
- [3] Paul Ekman. An argument for basic emotions. *Cognition and Emotion*, 6(3):169–200,1992.
- [4] S. Aman and S. Szpakowicz. . Using roget's thesaurus for fine-grained emotion recognition. In *International Joint Conference on Natural Language Processing*, 2008
- [5] Yanqing Chen and Steven Skiena. Building sentiment lexicons for all major languages. In *Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Short Papers)*, pages 383–389, Baltimore, Maryland, USA, June 23-25 2014. ACL, 2014.
- [6] Clematide S, Gindl S, Klenner M, Petrakis S, Remus R, Ruppenhofer J, Waltinger U, Wiegand M MLSA – a multi-layered reference corpus for german sentiment analysis, In: *Proceedings of the 8th international conference on language resources and evaluation (LREC'12)*, Istanbul, pp 3551–3556, 2012
- [7] Esuli A, Sebastiani F Determining the semantic orientation of terms through gloss classification. In: *Proceedings of the 14th ACM international conference on information and knowledge management, CIKM'05*, New York, pp 616–624, 2005
- [8] Virginia Francisco and Pablo Gervas. Automated mark up of aff ective information in english text. *Text, Speech and Dialouge*, volume 4188 of *Lecture Notes in Computer Science*:375–382, 2006.
- [9] David John, Anthony C. Boucouvalas, and Zhe Xu. Representing emotinal momentum within expressive internet communication. In *In Proceedings of the 24th IASTED international conference on Internet and multimedia systems and applications*, pages 183-188, Anaheim, CA, ACTA Press, 2006.
- [10] J. Johnson J. Guthrie K. Roberts, M.A. Roach and S.M. Harabagiu. 2012. "empatweet: Annotating and detecting emotions on twitter", In *in Proc. LREC, 2012*, pp.3806-3813.
- [11] Hatzivassiloglou V, McKeown KR Predicting the semantic orientation of adjectives. In: *Proceedings of the 35th annual meeting of the association for computational linguistics and 8<sup>th</sup> conference of the European chapter of the association for computational linguistics. Association for Computational Linguistics, Madrid*, pp 174–181, (1997)
- [12] Hu M, Liu B Mining opinion features in customer reviews. *AAAI* 4(4):755–760
- [13] Karttunen L, Zaenen A (2005) Veridicity. In: Katz G, Pustejovsky J, Schilder F (eds) *Annotating, extracting and reasoning about time and events. Internationales Begegnungsund Forschungszentrum für Informatik (IBFI), Schloss Dagstuhl*, 2004
- [14] Elsa Kim, Sam Gilbert, J.Edwards, and Erhardt Graeff . Detecting sadness in 140 characters: Sentiment analysis of mourning of michael jackson on twitter, 2009.
- [15] Kiritchenko S, Zhu X, Cherry C, Mohammad S Detecting aspects and sentiment in customer reviews. In: *Proceedings of the 8th international workshop on semantic evaluation exercises (SemEval-2014)*, Dublin, 2014
- [16] Klinger R, Cimiano P (2014) The USAGE review corpus for fine grained multi lingual opinion analysis. In: *Proceedings of the 9th International Conference on Language Resources and Evaluation*

(LREC'14), Reykjavik, pp 2009

[17] Saif M. Mohammad and Tony Yang. Tracking sentiment in mail : How genders differ on emotional axes. In Proceedings of the 2nd Workshop on Computational Approaches to Subjectivity and Sentiment Analysis(WASSA 2011), pages 70- 79, Portland, Oregon. Association for Computational Linguistics,2011.

[18] Saif Mohammad. #emotional tweets. In The First Joint Conference on Lexical and Computational Semantics – Volume 1: Proceedings of the main conference and the shared task, and Volume 2: Proceedings of the Sixth International Workshop on Semantic Evaluation (SemEval 2012), 2012.

[19] Alena Neviarouskaya, Helmut Prendinger, and Mitsuru Ishizuka. Recognition of affect, judgment, and appreciation in text. In Proceedings of the 23rd International Conference on Computational Linguistics, COLING '10, pages 806–814, Stroudsburg, PA, USA. Association for Computational Linguistics.2010.

[20] Pang B, Lee L A sentimental education: sentiment analysis using subjectivity based on minimum cuts. In: Proceedings of the 42nd meeting of the association for computational linguistics (ACL'04), Main Volume, Barcelona, pp 271–278 2003

[21] Pang B, Lee L, Vaithyanathan S Thumbs up? Sentiment classification using machine learning techniques. In: Proceedings of the 2002 conference on empirical methods in natural language processing (EMNLP 2002), Philadelphia, pp 79–86, 2002

[22] Lisa Pearl and Mark Steyvers. Identifying emotions, intentions and attitudes in text using a game with a purpose. In In Proceedings of the NAACLHLT 2010 Workshop on Computational Approaches to Analysis and Generation of Emotion in Text, Los Angeles, California.2010.

[23] Ashequl Qadir and Ellen Riloff . Bootstrapped learning of emotion hashtags #hashtags4you. In the 4th Workshop on Computational Approaches to Subjectivity, Sentiment & Social Media Analysis (WASSA 2013).2013.

[24] Remus R, Quasthoff U, Heyer G SentiWS – a publicly available German-language resource for sentiment analysis. In: Proceedings of the 7th international language resources and evaluation (LREC), Istanbul, 2010

[25] Carlo Strapparava and Alessandro Valitutti. Wordnet-affect: an affective extension of wordnet. Technical report, ITC-irst, Istituto per la Ricerca Scientifica e Tecnologica I-38050 PovoTrento Italy, 2004.

[26] Peter D. Turney Saif M. Mohammad. Crowdsourcing a word-emotion association lexicon. Computational Intelligence, 29 (3), 436-465, Wiley Blackwell Publishing Ltd, 2013, 29(3):436–465.2013.

[27] Taboada M, Brooke J, Tofiloski M, Voll K, Stede M Lexicon-based methods for sentiment analysis. Comput Ling 37(2):267–307, 2011

[28] Wenbo Wang, Lu Chen, Krishnaprasad Thirunarayan, and Amit P. Sheth. Harnessing twitter "big data" for automatic emotion identification. In Proceedings of the 2012 ASE/IEEE International Conference on Social Computing and 2012 ASE/IEEE.2012.

[29] C. Yang, K. H. Y. Lin, and H. H. Chen. Emotion classification using web blog corpora. In Proceedings of the IEEE/WIC/ACM International Conference on Web Intelligence, WI '07, pages 275–278, Washington, DC, USA. IEEE Computer Society, 2007.

[30] Radovan Garabik ir Indre Pileckyte L. Štúr Institute of Linguistics, Slovak Academy of Sciences, Bratislava, Slovakia,From Multilingual Dictionary to Lithuanian WordNet, 2013